

ÉCOLE NATIONALE DES CHARTES

---

Paul, Hector Kervegan

# Traitement, exploitation et analyse d'un corpus semi-structuré : le cas des catalogues de vente de manuscrits

Mémoire pour le diplôme de master

*Technologies numériques appliquées à l'histoire*

2022



# Résumé



# Introduction

J’ai choisi de structurer mon mémoire autour de plusieurs questions connexes, qui, à différents degrés, se retrouvent tout au long du développement :

En quoi la nature semi-structurée du corpus permet d’en automatiser le traitement ? Comment produire des informations normalisées et exploitables à partir d’un corpus textuel semi-structuré ? Comment rendre la recherche en humanités numériques réutilisable et encourager le partage de données entre projets de recherche ?

Les deux premières questions, d’orientation plutôt technique, forment la colonne vertébrale pour le mémoire ; elles lient deux aspects centraux : la nature du corpus et la manière dont sa structure permet toute la chaîne de traitement. Par « semi-structuré », j’entends que, à un niveau distant, toutes les entrées de catalogue suivent la même structure ; des séparateurs distinguent les différentes parties, et les informations sont souvent structurées de manière semblable pour chaque manuscrit vendu. Cela permet un traitement de « basse technologie » (*low-tech*) en évitant d’entraîner de lourds modèles de traitement du langage naturel (ce qui aboutirait à des solutions complexes, difficiles à maintenir et à faire évoluer et relativement opaques dans leur fonctionnement). À l’inverse, un corpus semi-structuré peut être traité en déduisant une « structure abstraite », que chaque entrée de catalogue partage. Il est alors possible de mettre en place des solutions techniques plus faciles, pour un résultat de qualité équivalente. Produire des « informations normalisées et exploitables » implique de traiter le corpus en cherchant des réponses à des questions de recherche précises – dans le cadre de mon stage, une question centrale a été de chercher à isoler les facteurs déterminant le prix d’un manuscrit.

Les deux dernières questions, au premier abord plus théoriques, me semblent centrales, notamment à la troisième partie de ce mémoire. Numérisation, traitement informatisé et diffusion sur le web ne sont pas des opérations neutres, mais un ensemble de « traductions » des documents originels. Ces processus comportent une part de choix conscients, qu’il s’agit de mettre en avant. Par exemple, on considère que la majorité des documents vendus ont pour titre l’auteur.ice du document. Cette personne n’est cependant pas toujours mentionnée, et des documents peuvent être nommés d’après un lieu, un événement ou un thème (la Révolution française, par exemple). Ces « traductions » des catalogues sont relativement discrètes tout au long de la chaîne de traitement (où le format dominant est la TEI, qui garde une relation d’équivalence avec le texte). C’est lors

du passage au site web que ce processus de traduction devient plus évident, et, potentiellement, plus problématique. On y abandonne la référence au document originel (les catalogues numérisés ne sont pas accessibles en ligne par un.e utilisateur.ice), le catalogue n'est plus la manière privilégiée d'accéder aux items vendus... De plus, la construction d'un site web implique la conception d'une interface et, dans notre cas, la production d'une série de visualisations intégrées au site. Le passage au site web remet aussi en cause la hiérarchie habituelle entre ingénierie et recherche : la conception d'un site ne répond pas à une question scientifique, mais elle soulève ses propres questions. Loin d'être anodines, ces problématiques de design déterminent la construction et la réception des savoirs. Il est donc important, je pense, de problématiser ces questions de visualisation et de design.

## Première partie

Du document numérisé au XML-TEI :  
nature du corpus, structure des  
documents et méthode de  
production des données





# Chapitre 1

## Le marché des manuscrits autographes au prisme des catalogues de vente

Ce chapitre présente l'objet d'étude du projet *MSS* : étudier le marché des manuscrits autographes du XIX<sup>ème</sup> s.. parisien à partir de ses catalogues de vente et étudier la construction du canon littéraire au prisme du marché du manuscrit.

### 1.1 Pourquoi étudier le marché des manuscrits autographes ?

Cette section porte sur l'intérêt scientifique des objets d'étude du projet (marché des manuscrits et étude de la construction du canon).

### 1.2 La structure du corpus : périodisation, producteurs des documents et classification

Ici est faite une présentation des documents traités dans le cadre du projet *MSS*. La présentation est à deux niveaux : au niveau du corpus et des catalogues. Ce chapitre s'appuie sur les mémoires effectués par d'ancien.ne.s stagiaires de *Katabase*, qui ont déjà beaucoup analysé la nature et les enjeux du corpus<sup>1</sup>.

---

1. Lucie Rondeau du Noyer, *Encoder automatiquement des catalogues en XML-TEI. Principes, évaluation et application à la revue des autographes de la librairie Charavay*, Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2019, URL : <https://github.com/lairaines/M2TNAH> (visité le 13/06/2022) ; Caroline Corbières, *Du catalogue au fichier TEI. Création d'un workflow pour encoder automatiquement en XML-TEI des catalogues d'exposition*. Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2020, URL : [https://github.com/carolinecorbieres/Memoire\\_TNAH](https://github.com/carolinecorbieres/Memoire_TNAH)

### 1.2.1 Le corpus de catalogues de vente de manuscrits

Ici est présenté le corpus : nature, quantité de documents (et d'entrées individuelles), dates, différentes classifications qui peuvent être faites (revues, catalogues de ventes aux enchères ou à prix fixes...).

### 1.2.2 Structure des catalogues

Ici sera présentée la structure des catalogues ; la structure de chaque page ne sera détaillée qu'à la partie suivante.

---

(visité le 13/06/2022) ; Juliette Janès, *Du catalogue papier au numérique. Une chaîne de traitement ouverte pour l'extraction d'information issue de documents structurés*, Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2021, URL : [https://github.com/Juliettejns/Memoire\\_TNAH](https://github.com/Juliettejns/Memoire_TNAH) (visité le 13/06/2022).

# Chapitre 2

## Production des données : de l’OCR à la TEI

Cette partie s’attache autant à présenter le processus d’océrisation (qui est déjà bien établi et ne constitue pas le cœur de mon stage) que la structure des documents. Alors que le chapitre précédent s’intéresse aux catalogues dans leur ensemble, ici, on étudie le corpus au niveau de la page et de l’entrée individuelle. En effet, l’océrisation repose sur la segmentation, et donc sur l’établissement d’une structure « abstraite » d’une page (c’est-à-dire, d’un découpage de la page en zones).

### 2.1 Extraire le texte des imprimés

#### 2.1.1 Comprendre la structure du document pour préparer l’édition numérique

##### Appréhender la structure de la page à l’aide de SegmOnto

La structure des catalogues est présentée au niveau de la page. L’ontologie SegmOnto<sup>1</sup> est utilisée, autant pour appréhender la structure de la page que pour exprimer cette structure de façon standardisée.

##### Description des entrées de catalogue : préparer l’édition TEI

Ici, la structure des catalogues est présentée au niveau de l’entrée, c’est à dire du lot mis en vente. C’est à partir de la structure des entrées qu’est construite l’édition XML-TEI. On s’intéresse à la structure des entrées individuelles à deux niveaux :

---

1. Kelly Christensen, Simon Gabay, Ariane Pinche et Jean-Baptiste Camps, « SegmOnto – A Controlled Vocabulary to Describe Historical Textual Sources », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris : École nationale des Chartes, 2022.

- Au niveau intellectuel : quelles sont les différentes parties d’une entrée (titre, description du manuscrit, prix...).
- Au niveau « textuel » : quels sont les séparateurs, c’est à dire les éléments dans le texte qui permettent de séparer les pages de catalogue en entrées et les entrées en sous-éléments) correspondant à la structure intellectuelle décrite ci-dessus.

## 2.2 L’encodage des manuscrits en XML-TEI

### 2.2.1 Encoder les catalogues en TEI

Ici est présentée la représentation XML-TEI des catalogues de vente.

### 2.2.2 L’encodage en TEI : un processus sélectif qui réduit les significations du texte

Après une étape d’océrisation via *eScriptorium*, le texte extrait des PDF peut être exporté soit en texte brut, soit en XML Page ou Alto. Ces formats s’attachent à garder une relation entre le XML et le document numérisé (les zones de texte sont indiquées, chaque ligne est dans une balise...). Cependant, l’unité intellectuelle centrale à la suite du projet, ce n’est pas la page numérisée, mais l’entrée de catalogue. Un format plus complexe que le XML d’*eScriptorium* est donc nécessaire. Assez logiquement, la suite du projet s’appuie sur une traduction des catalogues en TEI. On s’intéresse autant à la structure des documents XML (quelles balises sont utilisées...) qu’à l’intérêt scientifique d’une édition numérique (balisage sémantique, possibilité de normaliser les informations grâce à des attributs).

L’édition numérique en XML-TEI des catalogues implique une certaine perte d’informations : l’intégralité des significations contenues dans les catalogues imprimés ne peut être traduite en TEI (la police, ou la qualité du papier, peuvent être documentés mais ne peuvent pas être reproduites). Ce genre de perte d’information a lieu, à différents degrés, dans la plupart des éditions TEI : ce format n’est pas un substitut des documents originaux. Dans le projet *MSS / Katabase*, d’autres informations sont perdues : l’édition numérique n’est pas censée être une représentation exhaustive des catalogues. La TEI n’est pas utilisée comme un format de conservation, mais comme un format de traitement qui sera enrichi dans les différentes étapes. Afin de mesurer ce qui est conservé et ce qui est perdu du document original, l’édition TEI sera analysée à la lumière de la « roue du texte » du philologue Patrick Sahle<sup>2</sup> qui modélise les significations plurielles d’un texte.

cf Smith+Blackwell (2012) : texte (édition critique) peut être rpr sous la forme d’un graphe RDF de citations

---

2. Patrick Sahle, « Digital modelling. Modelling the digital edition », dans *Medieval and modern manuscript studies in the digital age*, London/Cambridge, 2016, URL : [https://dixit.uni-koeln.de/wp-content/uploads/2015/04/Camp1-Patrick\\_Sahle\\_-\\_Digital\\_Modelling.pdf](https://dixit.uni-koeln.de/wp-content/uploads/2015/04/Camp1-Patrick_Sahle_-_Digital_Modelling.pdf), p. 11.

parler aussi de l'OHCO (renear\_refining\_1996) et des critiques qui lui ont été faites dès les débuts de la TEI

John Frow + Susan Pearce ?



## Deuxième partie

Vers une étude des facteurs  
déterminant le prix des documents :  
alignement des entrées de catalogue  
et enrichissement sémantique avec  
*Wikidata*





Cette partie est construite à partir d’une question de recherche : comment produire des informations exploitables pour une étude économétrique à partir d’un corpus textuel semi-structuré ? Un des objectifs du projet est de faire l’étude des facteurs déterminant le prix d’un manuscrit. Pour faire cette étude, il faut obtenir, pour chaque entrée du catalogue, des données normalisées et exploitables. Le travail d’extraction de données présentes dans les catalogues a déjà été fait par de précédent.e.s stagiaires. Ces données sont principalement quantitatives : prix des manuscrits, dimensions et nombre de pages, date de création. Il est nécessaire de compléter les informations par des données qualitatives et d’enrichir les catalogues sources extérieures : c’est, enrichissement, qualifié de sémantique, qui permettra d’étudier la relation entre facteurs biographiques et prix d’un catalogue. Pour réaliser cet enrichissement, il a été choisi d’associer les auteur.ice.s des manuscrits avec des entités dans une base de connaissance en ligne : *Wikidata*. Concrètement, dans une base de connaissance, une entité est représentée par un identifiant auquel sont liés d’autres données. Il s’agit donc de récupérer des identifiants depuis *Wikidata* pour récupérer automatiquement des informations sur les personnes via *SPARQL*. Le choix de travailler uniquement sur les noms, et non sur la description des documents, a deux motivations :

- Les noms de personnes (et la manière dont elles sont décrites) constituent la partie la plus normalisée des documents. La description des manuscrits est plutôt en « texte libre ». Dans la continuité avec le reste du projet, une approche de « basse technologie », qui consiste à s’appuyer majoritairement sur des solutions techniquement simples, a été privilégiée.
- Toutes les informations « simples » (données quantitatives facilement normalisables : dates etc.) ont déjà été extraites des descriptions des manuscrits.

Ce travail d’enrichissement a été fait en deux temps. La première étape, et la plus difficile, est l’alignement avec *Wikidata*. Cela demande d’extraire un ensemble d’informations à partir du nom de la personne et de la description de celle-ci. Parmi les informations extraites : nom, prénom, titre de noblesse, occupation, dates de vie et de mort. À partir de ces informations, un algorithme lance plusieurs recherches sur *Wikidata* afin d’aligner les données des catalogues avec les bonnes entités. La deuxième étape, nettement plus simple, consiste à lancer des requêtes *Wikidata* sur les identifiants récupérés afin de récupérer des informations sur les auteur.ice.s des manuscrits (cette partie du travail est encore en cours) pour enrichir nos données. Pour finir, les fichiers TEI sont mis à jour pour ajouter les identifiants *Wikidata*. Ainsi, il est possible de faire le lien entre les entrées de catalogues dans des fichiers XML et les données issues de requêtes *SPARQL*, stockées dans un JSON.

Mais au delà de la question de recherche présentée en début de partie, et des problématiques techniques qui seront détaillées plus bas, cette partie s’intéresse également aux méthodes utilisées en humanités numériques, et cherche à comparer la méthode établie ici

avec d'autres techniques plus répandues, mais aussi plus complexes et basées sur l'apprentissage machine, pour lier des textes avec des bases de connaissance. En comparant les méthodes utilisées, la possibilité de les adapter à d'autres corpus et les résultats obtenus, cette partie s'attache à étudier la possibilité de développer des méthodes « alternatives », basées sur la structure des documents traités, pour le traitement de données textuelles, tout en montrant comment les méthodes établies peuvent servir d'inspiration pour les possibilités offertes par l'enrichissement sémantique.

# Chapitre 3

## Questions introductives : pourquoi et comment s'aligner avec *Wikidata* ?

Cette section, introductive, répond à des questions évidentes mais essentielles : elles permettent de mettre au clair l'intérêt et les (multiples) difficultés dans l'alignement avec *Wikidata*.

### 3.1 Questions théoriques

#### 3.1.1 Pourquoi s'aligner avec des entités externes ?

L'enrichissement des données depuis une base de connaissances a pour objectif principal de mieux comprendre les déterminants du prix des manuscrits sur le marché du XIX<sup>ème</sup> s.. Mais pourquoi passer par cet alignement ?

Pour étudier les déterminants du prix d'un manuscrit, il faut établir la relation entre le prix d'un manuscrit et un ensemble d'autres facteurs (qui a écrit un manuscrit, quelles sont ses dimensions, de quand date le document...). En d'autres termes, il faut étudier le comportement d'une variable en fonction d'autres variables. En économétrie, cette opération s'appelle le calcul de régressions linéaires. La variable étudiée (le prix) est dite la variable expliquée ; les facteurs déterminant sa valeur sont dites « variables explicatives »<sup>1</sup>. Cependant, cette opération est loin d'être anodine : il faut d'abord identifier les variables pertinentes, et ensuite trouver un moyen de les quantifier. Deux difficultés se présentent pour alors.

Premièrement, il faut pouvoir quantifier les variables expliquées pour calculer des régressions linéaires. Il est possible de leur assigner une valeur numéraire (ce qui est aisé pour les informations quantitatives des catalogues : la date de l'écriture d'un manuscrit, ses les dimensions). Une autre possibilité est de quantifier la présence ou non d'une variable :

---

1. *Régression linéaire*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/R%C3%A9gression\\_lin%C3%A9aire](https://fr.wikipedia.org/wiki/R%C3%A9gression_lin%C3%A9aire) (visité le 10/07/2022).

mention d'un.e destinataire ou du contenu d'un manuscrit. Cependant, ces approches quantitatives ne permettent pas de quantifier des informations complexes, comme la célébrité des auteur.ice.s, ou encore si un manuscrit porte sur un évènement historique ou biographique important (le manuscrit d'un texte célèbre, par exemple, pourrait avoir une valeur particulières). Ces informations sont parfois être présentes dans les catalogues ; elles peuvent aussi être connues des lecteur.ice.s d'aujourd'hui et des acheteurs et acheteuses de l'époque. Il n'existe cependant pas de méthodes faciles pour détecter ou quantifier la célébrité d'une personne, ou l'importance d'un sujet.

Une deuxième difficulté découle justement de la part d'implicite qu'il y a dans les catalogues. Les descriptions des items vendus sont brèves, et comprendre ce qui fait la valeur d'un manuscrit demande aux acheteur.euse.s d'avoir des références culturelles et historiques : celles-ci permettent d'identifier l'auteur.ice ou le sujet, et donc pour comprendre la valeur d'un manuscrit. Dans le cadre du projet *MSS / Katabase*, les entrées de catalogues sont traitées par une machine qui, en toute logique, ne dispose pas de ces références. La compréhension qualitative des entrées de catalogues n'est donc pas compatible avec l'approche par lecture distante qui permet l'analyse de très larges corpus textuels. Pour éviter de perdre ces informations qualitatives essentielles, il est donc nécessaire de trouver un moyen de quantifier le qualitatif.

En bref, la question est : comment faire la différence entre une lettre de La Rochefoucauld (3.1a), vendue 200 francs, et la deuxième (3.1b), vendue à 30 francs ? Le problème est un problème de lecture. Une observation de la description des lettres par un être humain comme par une machine peuvent identifier des éléments semblables dans le texte : les deux lettres sont écrites par des ducs ; l'une est une « Très-belle lettre » (3.1a), l'autre est une « Lettre intéressante » (3.1b). Bien que les lettres partagent des attributs, il y a une forte différence de prix entre les deux manuscrits. Un.e lecteur.ice peut trouver une raison à cette différence de prix : La Rochefoucauld et Madeleine de Scudéry n'ont pas le même statut que le duc de Villars. Un regard humain peut donc interpréter un prix et déterminer une valeur en s'appuyant sur ses connaissances. La lecture est qualitative et s'appuie sur de l'implicite, ce qui n'est pas possible pour une machine : formellement, rien ne distingue un nom d'un autre ; lorsqu'un programme « lit » un texte, il ne peut pas s'appuyer sur ses connaissances pour déterminer ce qu'un nom signifie, ce à quoi il fait référence.

Pour analyser efficacement la variable « prix », il faut pourtant pouvoir, dans une certaine mesure, comprendre les informations implicites et qualitatives contenues dans les catalogues. Le parti pris a donc été de construire le socle de connaissance qui manque à une machine, en s'alignant avec *Wikidata* et en s'en servant pour enrichir nos données. Le choix a été fait de ne s'aligner avec *Wikidata* que pour certaines parties des entrées de catalogue. Pour rappel, voici leur structure (3.1) :

Les entrées de catalogue contiennent beaucoup d'informations qualitatives, qui pour-

```

1 <item n="287" xml:id="CAT_000126_e287">
2   <num type="lot">287</num>
3   <name type="author">Tascher de la Pagerie
4     ↪ (Marie-Euphémie-Désirée)</name>
5   <trait>
6     <p>fille de Joseph Tascher qui passa à Saint-Domingue, et de
7       ↪ Mlle de la Chevalerie; elle épousa M. de Renaudin, puis le
8       ↪ marquis François de Beauharnais et contribua beaucoup à la
9       ↪ fortune de la future impératrice Joséphine.</p>
10  </trait>
11  <desc xml:id="CAT_000126_e287_d1">
12    Pièce originale;
13    (<date when="1780">1780</date>),
14    <measure type="length" unit="p" n="4">4 p.</measure>
15    <measure type="format" unit="f"
16      ↪ ana="#document_format_4">in-4.</measure>
17  </desc>
18  <measure commodity="currency" unit="FRF" quantity="15">15</measure>
19  <note>Plaidoirie pour le dit Renaudin contre sa femme qu'on accusait
20    ↪ d'avoir usé de l'influence de son parent, M. de Beauharnais,
21    ↪ pour se faire épouser</note>
22 </item>

```

Code source 3.1 – Représentation XML-TEI d'une entrée de catalogue

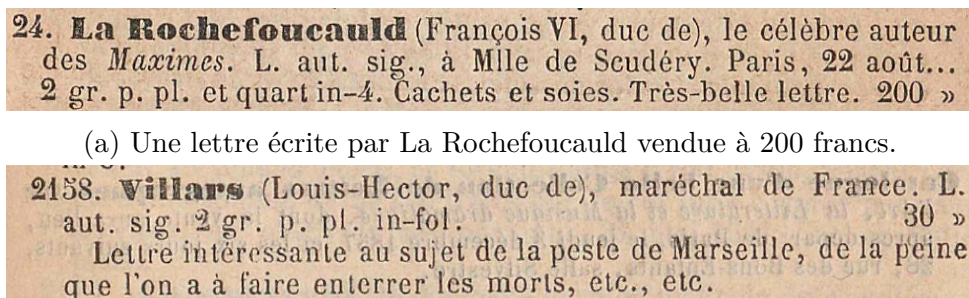


FIGURE 3.1 – Deux exemples de lettres

raient avoir une influence sur le prix du manuscrit : ici par exemple, la description du contenu de la lettre dans le `tei:note` ; il est également souvent fait mention du ou de la destinataire. Cependant, l'alignement avec *Wikidata* n'a pas été fait avec l'intégralité des entrées, mais seulement le contenu du `tei:name`. Le `tei:desc` a déjà fait l'objet d'un grand travail de normalisation et d'extraction d'informations ; un alignement avec des sources externes n'aurait donc pas une très grande plus-value. L'élément `tei:note` contient souvent des informations intéressantes, puisque c'est là qu'est décrit le contenu d'un manuscrit. Cependant, cet élément n'est pas toujours présent ; son contenu est souvent écrit en langage naturel, non structuré, et contient des informations trop variées pour développer un traitement uniforme. Il est donc difficile de tirer parti de cet élément. Le `tei:name` et son `tei:trait` sont les éléments les plus régulièrement présents ; les informations qu'ils contiennent sont toujours les mêmes (nom d'une personne ou thème d'un manuscrit dans le `tei:name`, description du `tei:name` dans le `tei:trait`) ; enfin, ces deux éléments n'ont pas du tout été transformés dans le reste de la chaîne de traitement. Ils portent donc des informations qualitatives centrales pour produire des données exploitables dans une étude économétrique.

Le parti pris a donc été d'aligner avec des identifiants *Wikidata* les noms contenus dans les balises `tei:name` à l'aide des descriptions contenues dans les `tei:trait` ; à partir de cet alignement a été constituée une base de données. Cela permet d'approximer une lecture « humaine » des items en vente : pour chaque auteur.ice, un certain nombre d'informations auront été récupérées pour mieux identifier la personne (ses occupations, son origine, ses dates de vie...). L'analyse du corpus s'appuie alors sur un bagage de connaissances qui permet d'appréhender par lecture distante l'importance d'une personne. Il devient alors envisageable de voir dans quelle mesure la mention d'une personne impacte le prix d'un manuscrit, et quels sont les facteurs biographiques déterminant dans l'établissement de la valeur. Pour revenir à l'exemple de La Rochefoucauld : à défaut de permettre de savoir qui il est, un alignement avec *Wikidata* permet de mieux situer son statut et sa place dans la culture française, en récupérant le nombre de ses publications ou encore les institutions dont il est membre.

### 3.1.2 Pourquoi utiliser la base de connaissances *Wikidata* ?

De nombreuses bases de connaissances sont accessibles en ligne ; celles-ci ont pour rôle de stocker et de diffuser des données structurées. Dans le paysage de la diffusion d'information sur le Web, les bases de connaissance s'opposent donc à de nombreux projets qui diffusent du texte, comme l'encyclopédie *Wikidata*. Elles s'opposent également aux « dumps » de données, c'est-à-dire à la diffusion de jeux de données publics sur le Web. Les « dumps » – tels qu'ils sont pratiqués par *Data.gouv*, structure développée par la France – sont des jeux de données « inertes » : les données sont diffusées dans des fichiers qui, une fois diffusés, ne seront pas mis à jour. Ces bases de connaissances donnent accès à de l'information brute, directement exploitable ; les informations y sont accessibles individuellement et directement sur le Web (alors qu'elles sont contenues dans de lourds fichiers qui doivent être téléchargés dans le cas des « dumps »). Certaines de ces bases de connaissances ont été créées par des institutions publiques, comme *DataBnF*, créée par la Bibliothèque nationale de France. D'autres ont des sources universitaires, comme *DBPedia* ou communautaires, c'est-à-dire qu'elles sont enrichies par une communauté d'utilisateur.ice.s – dont certains peuvent être des acteurs institutionnels. Enfin, certaines de ces bases de connaissances sont spécialistes, comme le Thésaurus INHA, qui met à disposition du public les vocabulaires contrôlés développés par l'Institut d'histoire de l'art ; d'autres, comme *DBPedia*, sont généralistes. Un paysage varié – brièvement présenté ici – de solutions existent donc pour récupérer des données. Parmi cette variété de sources, laquelle choisir ? Cette question fait l'objet de nombreuses études, notamment issues du domaine du traitement automatisé du langage<sup>2</sup>. Des algorithmes complexes ont été développés pour juger de la qualité des données et de leur adaptation aux corpus d'un projet de recherche. L'utilisation de tels algorithmes ne s'intègre pas dans les méthodes du projet *Katabase* : ces algorithmes sont faits pour ensuite aligner les documents sources avec une base de connaissance à l'aide de méthodes d'apprentissage machine. Le projet *MSS / Katabase* évite, lorsque cela est possible d'utiliser de tels outils. Une chaîne de traitement pour l'alignement avec *Wikidata* a de plus déjà été élaborée ; il n'était pas pertinent de modifier toute la chaîne de traitement simplement pour déterminer si une meilleure base de connaissance pourrait être utilisée.

La sélection de la base de connaissance à utiliser a donc été faite à partir d'une évaluation des besoins du projet. D'un point de vue des données, il était nécessaire d'utiliser une base généraliste, disposant d'informations biographiques sur des personnes et de données pertinentes pour un corpus issu de la France métropolitaine du XIX<sup>ème</sup> s.. D'un point de vue technique, la base de connaissance devait disposer de deux outils. D'abord, l'alignement des `tei:name` avec des entités *Wikidata* devait être fait à l'aide de recherches

---

2. Carmen Brando, Nathalie Abadie et Fontini Frontini, « Évaluation de la qualité des sources du web de données pour la résolution d'entités nommées », *Ingénierie des systèmes d'information*, 21–5 (28 déc. 2016), p. 31–54, DOI : 10.3166/isi.21.5–6.31–54.

en plein texte sur le moteur de recherche de la base ; pour automatiser ce processus, la base de connaissance devait donc disposer d’une Interface de programmation d’application (API) qui permette une telle opération. Ensuite, la constitution d’une base de donnée sur ces entités devait être faite à l’aide du langage de requêtes **SPARQL**, qui permet la récupération automatisée d’informations structurées sur une base de données en ligne. *Wikidata* répond à ces deux conditions. Mais cette base de donnée communautaire, développée par la *Wikimedia Foundation* (également à l’origine de *Wikipédia*) et rendue publique en 2012 présente de nombreux autres avantages. Contrairement à *DBPedia*<sup>3</sup>, les données de *Wikidata* sont directement produites par les contributeur.ice.s. Le projet est bicéphale, composé d’un côté de l’ajout d’informations à destination du public, et de l’autre de la construction d’une ontologie ; la particularité de *Wikidata* est que les deux aspects sont collaboratifs, alors qu’en général, la construction d’une ontologie est le fait d’ingénieur.e.s professionnel.le.s<sup>4</sup>. En 2015, la communauté *Wikidata* représente 300000 personnes et 160 *bots* (des programmes informatiques qui ajoutent automatiquement des données sur *Wikidata*)<sup>5</sup>. Cette communauté est basée sur une forme de division du travail entre humain.e.s et machines : 85% des ajouts sont faits par des *bots*, mais ce sont des ajouts simples et automatisables (par exemple, l’ajout d’identifiants uniques pour lier les entités à d’autres sources de données). La part humaine des contributions est consacrée aux ajouts plus complexes, avec une minorité d’acteur.ice.s (30%) qui réalisent la très grande majorité des éditions, et sur de nombreux projets. À l’inverse, Les 70% restants ont une contribution bien plus faible, et centrée sur un ou quelques sujets de prédilection. L’immense majorité des éditions représente des ajouts d’informations, tandis que l’évolution de la structure de *Wikidata* repose sur une poignée d’utilisateur.ice.s très spécialisé.e.s<sup>6</sup>. Rien que par son envergure et la quantité de données disponibles, *Wikidata* est un candidat intéressant pour le projet d’alignement des **tei:name** avec une base de connaissances externe. Par ailleurs, sa nature généraliste s’adapte à la vaste période couverte par le corpus, ainsi qu’à la variété d’informations qui peuvent être contenues dans le **tei:name**. La nature communautaire du projet est elle aussi intéressante : *Wikidata* est enrichi par une communauté internationale, et de nombreuses données sont donc disponibles pour un corpus français ; des bases non-communautaires risquent toujours d’avoir un biais national dans les données disponibles.

Le modèle de données et les choix techniques pour leur ingénierie font de *Wikidata* une source adaptée aux objectifs du projet *Katabase*. Tous les objets qui sont documentés

---

3. Claudia Müller-Birn, Benjamin Karran, Janette Lehmann et Markus Luczak-Rösch, « Peer-production system or collaborative ontology engineering effort : what is Wikidata ? », dans *Proceedings of the 11th International Symposium on Open Collaboration*, San Francisco California, 2015, p. 1-10, DOI : 10.1145/2788993.2789836, p. 1.

4. *Ibid.*, p. 2-3.

5. *Ibid.*, p. 4.

6. *Ibid.*, p. 4, 8-9.



par *Wikidata* (appelés entités) peuvent être classés en deux types<sup>7</sup> :

- Les items, c'est-à-dire des données sur laquelle d'autres informations sont disponibles. Ceux-ci correspondent à des choses physiques ou à des concepts du monde réel, et peuvent être classifiés en deux catégories : les classes et les individus. Les classes sont des catégories génériques (« une personne ») tandis que les individus sont des membres de ces catégories (« Natalia Gontcharova »).
- Les propriétés, qui caractérisent les relations entre un ou plusieurs items. Par exemple, « occupation » permettrait de lier l'item « Natalia Gontcharova » à l'item « peintre ».

D'un point de vue technique, ce modèle est intéressant : il correspond à l'organisation des données du Web sémantique et de SPARQL<sup>8</sup> ; cela facilite l'utilisation ultérieure de ce langage de requêtes (décrite plus bas) : la relation entre les données, qui est utilisée pour récupérer des informations avec SPARQL, est explicite ; il est donc facile de trouver comment structurer des requêtes. Là encore en concordance avec les modèles de données sémantiques, chaque item dispose d'un identifiant unique, d'un *Uniform Resource Identifier* (URI) ; c'est ces identifiants qui sont utilisés dans SPARQL pour récupérer des données, ce qui facilite encore l'utilisation de ce langage<sup>9</sup>. Enfin, une dernière spécificité de cette base de connaissances est intéressante en vue de l'utilisation ultérieure de SPARQL : *Wikidata* a fait le choix de n'utiliser qu'un vocabulaire par catégorie de données. Comme nous le verrons plus tard, dans le Web sémantique, un vocabulaire est nécessaire pour définir les propriétés reliant plusieurs entités ; par exemple, un identifiant est relié à son nom grâce à la propriété « `rdfs:label` », qui appartient au vocabulaire `rdfs`. Ces vocabulaires doivent être utilisés avec SPARQL. Plusieurs vocabulaires ont souvent leurs propres manières de définir les mêmes relations, et les bases de connaissance ont tendance à utiliser plusieurs vocabulaires qui se chevauchent – c'est par exemple le cas de *DataBnF*. Cela rend le modèle de données utilisés par ces bases particulièrement complexe ; il est donc difficile de naviguer dedans, et encore plus difficile de récupérer des données via SPARQL. À l'inverse, la majorité des propriétés sont définies par l'ontologie de *Wikidata*, qui est très bien documenté sur le site ; cela rend la manipulation des données sur cette base de connaissances particulièrement aisée. Ainsi, *Wikidata* dispose d'un modèle de donnée totalement adapté à l'utilisation de langages de requêtes dédiées aux modèles de données sémantiques.

Dans *Wikidata*, tout item contient les mêmes types informations. Trois sont particulièrement intéressantes pour le processus d'alignement : le label, la description et les

---

7. Fredo Erxleben, Michael Günther, Markus Krötzsch, Julian Mendez et Denny Vrandečić, « Introducing Wikidata to the Linked Data Web », dans *The Semantic Web – ISWC 2014*, dir. Peter Mika, Tania Tudorache, Abraham Bernstein, Chris Welty, Craig Knoblock, Denny Vrandečić, Paul Groth, Natasha Noy, Krzysztof Janowicz et Carole Goble, Series Title : Lecture Notes in Computer Science, Cham, 2014, t. 8796, p. 50-65, DOI : 10.1007/978-3-319-11964-9\_4, p. 51-52.

8. *Ibid.*

9. *Ibid.*, p. 56.

alias<sup>10</sup>. Un label permet de nommer un item ; par conséquent, il est possible de distinguer l'entité (représentée par un identifiant normalisé) et sa représentation linguistique (le label). Cette distance est intéressante pour le projet *Katabase* : l'alignement des `tei:name` avec des entités *Wikidata* revient retrouver un identifiant à partir du label qui lui est associé. Ce label est unique ; cependant, il existe de nombreux alias, qui sont des représentations linguistiques multilingues ou des variations du nom associé à une entité. Cela veut dire que, même si, dans un catalogue, un.e auteur.ice n'est pas nommé.e de façon canonique (en traduisant son nom, par exemple), il est toujours possible d'aligner le `tei:name` avec une entité *Wikidata* à l'aide des alias. La description, enfin, contient souvent la date de naissance et de décès ainsi que l'occupation d'une personne ; autant d'informations qui sont contenues dans le `tei:trait` et qui seront donc utilisées dans le processus d'alignement. De l'expérience que j'ai pu en faire, label, alias et description sont indexés par le moteur de recherche de *Wikidata* : cela veut dire que en recherchant ces informations sur l'API, les données des catalogues pourront être alignées avec les entités pertinentes.

Pour toutes ces raisons, la base de connaissance pour le processus d'alignement avec une source de données externe est *Wikidata*. Cette plateforme dispose à la fois des outils techniques nécessaires, d'un socle de données correspondant à l'univers de discours du corpus et d'une certaine facilité d'utilisation, un critère non négligeable pour garantir la réutilisation ultérieure des processus développés.

### 3.1.3 Quelle relation avec la résolution d'entités nommées d'entités nommées ?

En traitement automatisé du langage, le processus d'alignement d'éléments d'un texte avec une base de connaissances externe a un nom : c'est la résolution d'entités nommées (plus tard abrégée sous son acronyme anglo-saxon, *Named Entity Linking* (NEL)). Pour reprendre la définition proposée par Brando et al. : « La tâche dite de résolution d'entités nommées [...] consiste à affecter automatiquement la bonne référence, l'identifiant de la ressource correspondante dans une base de connaissances ou bien NIL (l'absence de valeur), à une entité nommée préalablement identifiée dans un texte et étiquetée. »<sup>11</sup>. Dans un texte, ce processus arrive après l'identification dans un texte d'entités, c'est-à-dire d'objets uniques et signifiants pour une problématique de recherche : personnes, lieux, mais aussi parfois concepts. Il consiste à aligner les entités du texte avec celles présentes sur une base de connaissance externe. Deux objectifs primaires peuvent être identifiés<sup>12</sup>. Le premier est la désambiguïsation des entités dans un texte, ce qui permet

---

10. *Ibid.*, p. 52.

11. C. Brando, N. Abadie et Fontini Frontini, « Évaluation de la qualité des sources du web de données pour la résolution d'entités nommées »..., p. 31.

12. *Ibid.*, p. 32.

à une machine de ne pas confondre deux entités homonymes (Alexandre Dumas père et fils, par exemple). Cette étape est très importante pour le traitement automatisé d'un texte : si un.e lecteur.ice humain.e est capable de distinguer les deux Alexandre Dumas en fonction de sa lecture et de ses références, cette distinction est impossible par une machine. La reconnaissance d'entités nommées permet alors d'explicitier l'implicite, ce qui est également l'un des objectifs ayant mené à l'alignement des auteur.ice.s de catalogues avec des entités de *Wikidata*. Le deuxième objectif est « l'enrichissement sémantique », c'est-à-dire la récupération d'informations depuis une source externe pour permettre une meilleure compréhension et analyse du texte par des méthodes computationnelles. Comme cela a été dit plus haut, c'est exactement à cette fin qu'est mené l'alignement des `tei:name` avec *Wikidata*. Pourquoi alors ne pas directement parler d'un processus de résolution d'entités nommées ?

Tout d'abord, l'alignement avec *Wikidata* n'était pas censé prendre une telle ampleur à ses débuts. Il avait initialement été prévu de réaliser cet alignement « à la main », ce qui, pour un corpus de près de 100000 entrées, aurait été extrêmement chronophage. C'est avant tout pour répondre à ce problème purement pratique qu'ont été développées des méthodes pour automatiser l'alignement qui sont présentées ici. Étant donné que les entrées sont riches d'informations, il est apparu que ce processus pouvait être réalisé automatiquement, avec bien sûr un certain degré d'erreur. En choisissant d'automatiser l'alignement, les méthodes d'extraction d'informations des catalogues et de recherches sur l'API se sont progressivement complexifiées afin de mieux s'adapter au corpus. L'alignement avec *Wikidata* a donc progressivement été développé, jusqu'à devenir un processus analogue à la résolution d'entités nommées.

Ensuite et surtout, la résolution d'entités nommées « traditionnelle » implique des méthodes en tout point différentes de celles qui sont présentées ici. Le NEL repose très fortement sur l'apprentissage machine. Il est d'abord utilisé pour identifier les entités dans le texte<sup>13</sup> ; ensuite, il sert à mener l'alignement des entités avec la base de connaissances elle-même. Pour ce processus, de nombreux outils ont été développés, comme REDEN<sup>14</sup> ou DBpediaSpotlight<sup>15</sup>, conçu au sein du projet *DBpedia* pour la résolution d'entités en utilisant cette base de connaissance. Ces outils fonctionnent, d'une certaine manière, à l'inverse du processus mis en place ici : à partir d'un ensemble de références potentielles récupérées d'une base de connaissance à l'aide de requêtes SPARQL<sup>16</sup> avant d'extraire une

---

13. *Ibid.*

14. C. Brando, Francesca Frontini et Jean-Gabriel Ganascia, « REDEN : Named Entity Linking in Digital Literary Editions Using Linked Data Sets », *Complex Systems Informatics and Modeling Quarterly* 7 (29 juil. 2016), p. 60-80, DOI : 10.7250/csimq.2016-7.04.

15. Pablo N. Mendes, Max Jakob, Andrés García-Silva et Christian Bizer, « DBpedia spotlight : shedding light on the web of documents », dans *Proceedings of the 7th International Conference on Semantic Systems - I-Semantics '11*, Graz, Austria, 2011, p. 1-8, DOI : 10.1145/2063518.2063519.

16. C. Brando, N. Abadie et Fontini Frontini, « Évaluation de la qualité des sources du web de données pour la résolution d'entités nommées »..., p. 35.

liste de candidats potentiels pour la résolution d'entités, puis de résoudre l'entité en y assignant un item issu de la base de connaissance<sup>17</sup>. À l'inverse, le processus mis en place ici consiste d'abord à extraire des données des textes, ensuite d'interagir avec un moteur de recherche pour sélectionner une entité parmi toutes celles disponibles sur la base de connaissances, et enfin de mener des requêtes SPARQL pour récupérer des données supplémentaires sur les items liés. Dans le détail, les méthodes utilisées pour la résolution d'entités sont également très différentes : elles se basent, de manière générale, sur une tentative d'aligner l'entité du texte avec une entité sur la base de connaissance par l'analyse de graphes<sup>18</sup>. Pour ce faire, les entités de la base sont représentées sous forme de graphes qui sont analysés pour trouver, au sein du graphe, l'entité pertinente ; l'apprentissage machine permet à un outil « d'apprendre » à reproduire tout seul ce processus d'analyse du graphe, afin d'aligner d'autres entités avec leur équivalent sur une base de connaissances. Dans le processus décrit ici, il n'y a ni analyse de graphes, ni apprentissage machine : l'alignement se base sur la détection de motifs à l'aide d'expressions régulières. Des informations sont détectées dans les catalogues, extraites et si besoin retraitées, avant d'être utilisées pour faire des recherches sur *Wikidata* afin d'aligner une auteur.ice avec une entité de cette base de connaissances. Le processus décrit ici privilégie donc une compréhension approfondie du corpus, de sa structure et de son vocabulaire, ainsi qu'une bonne connaissance du fonctionnement du moteur de recherche de *Wikidata*. Si les méthodes sont radicalement différentes, les possibilités ouvertes sont équivalentes. Les scores de performance des deux méthodes sont, comme nous le verrons, relativement proches, ce qui ouvre à la possibilité de chercher à développer des méthodes alternatives pour la résolution d'entités nommées.

## 3.2 Questions de méthode : comment faire le lien entre un corpus historique et une base de connaissances en ligne ?

### 3.2.1 Présentation générale de l'algorithme

Construire un jeu de données issu de SPARQL à partir de la manière dont une personne est nommée et décrite au XIX<sup>ème</sup> s. n'est pas une opération anodine. La chaîne de traitement est donc assez complexe, comme le montre le schéma 3.2. Cette chaîne de traitement peut être séparée en trois étapes.

---

17. *Ibid.*, p. 33.

18. C. Brando, Francesca Frontini et J.G. Ganascia, « REDEN... », p. 63.

### Étape 1 – Extraction et structuration de données

Premièrement, il s’agit d’aligner les entrées de catalogue avec des identifiants *Wikidata*. Ceux-ci sont liés à des entités *Wikidata*. Cette première étape repose avant tout sur l’extraction et la traduction des données depuis les éléments `tei:name` et `tei:trait`. Ce processus d’extraction permet de récupérer toutes les données pertinentes pour chaque entrée de catalogue et de les stocker dans un dictionnaire structuré. Comme on le verra, la nature « semi-structurée » des entrées (ainsi qu’une bonne connaissance du corpus) permet de d’automatiser le processus d’extraction et de traduction des données par détection de motifs, sans avoir à passer par l’apprentissage machine : étant donné que les mêmes types d’informations sont toujours présentes et que les entrées suivent des modèles relativement proches, il est possible de s’appuyer sur la structure des entrées pour identifier les informations pertinentes. L’extraction de données repose donc sur de la détection de motifs à l’aide d’expressions régulières : des récurrences sont repérées dans le texte et utilisées pour distinguer différentes informations (nom, prénom, titre de noblesse...). Pour appuyer l’usage d’expressions régulières par une méthode plus « qualitative » et précise, certains termes particuliers sont extraits et éventuellement traduits à l’aide de tables de conversion (c’est-à-dire de dictionnaires qui associent à un terme dans le texte une version normalisée).

### Étape 2 – récupération d’identifiants *Wikidata* via des recherches en plein texte à l’aide d’une API

Une fois les données du `tei:name` et du `tei:trait` structurées en dictionnaire, elles sont utilisées pour lancer plusieurs recherches en plein texte sur le moteur de recherche de *Wikidata*. Ces recherches sont faites automatiquement grâce à l’API de *Wikidata*. Pour maximiser les chances d’obtenir un identifiant valide, un algorithme a été conçu pour lancer plusieurs recherches à partir de chaque entrée. La première recherche met bout-à-bout toutes les valeurs disponibles dans le dictionnaire. Ensuite, en fonction des paramètres de recherche disponibles dans le dictionnaire, différentes autres recherches sont lancées. Cet algorithme a été élaboré en menant de nombreux tests pour maximiser le taux de réussite, calculé sous la forme d’un score F1.

### Étape 3 – constitution d’un jeu de données à l’aide de SPARQL

Si l’étape d’alignement avec *Wikidata* est la plus complexe, elle n’est qu’une étape préparatoire vers la constitution du jeu de données et à l’enrichissement sémantique. En fait, résoudre les entités est ce qui rend possible l’enrichissement en tant que tel : en lançant une requête SPARQL sur tous les identifiants obtenus, il est possible de récupérer des informations depuis *Wikidata* et donc de construire le jeu de données définitif. Pour cette étape, le processus est plus simple : les identifiants récupérés à la fin de l’étape

précédente sont dédoublonnés (pour éviter de lancer plusieurs fois la même requête) ; ensuite une requête **SPARQL** est initialisée et lancée chacun des identifiants. Les résultats sont traduits depuis les formats **JSON** ou **XML** retournés par **SPARQL** sous forme de **JSON** plus simple, et donc plus aisément manipulable. Le jeu de données est enregistré dans un fichier. Pour finir, les identifiants *Wikidata* sont réinjectés aux catalogues **TEI**, afin de pouvoir faire le lien entre les documents et le jeu de données qui a été construit.

Cette chaîne de traitement étant lancée sur plus de 80000 entrées de catalogues, le temps d'exécution est très long et même des petites améliorations de performance peuvent avoir un grand impact ; dans sa version initiale, le script demandait des performances particulièrement élevées, et ne fonctionnait pas sur un ordinateur de bureau. La chaîne de traitement a donc été reprise en plusieurs points afin d'être optimisée, de fonctionner plus vite en étant moins coûteuse en ressources.

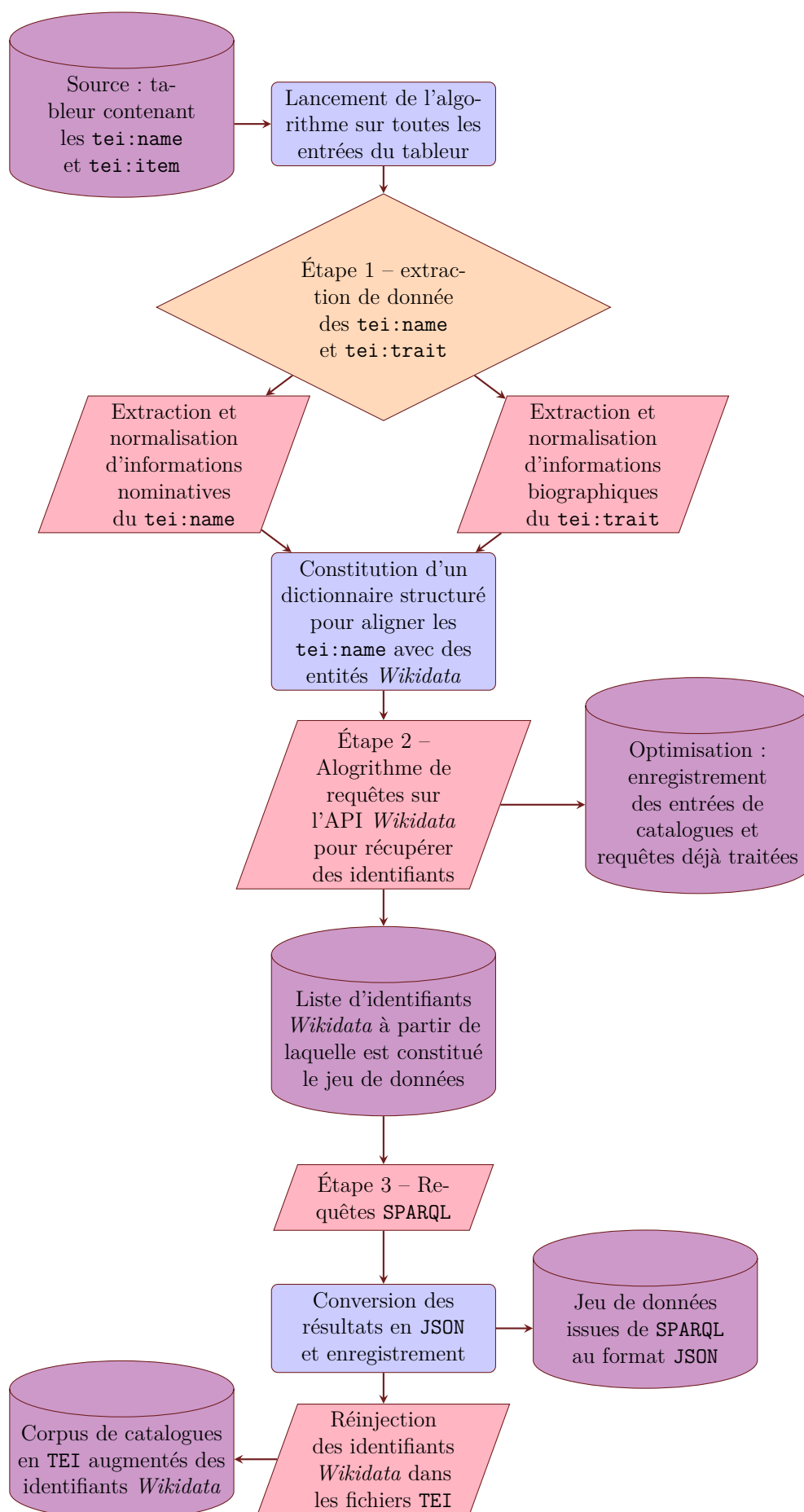


FIGURE 3.2 –  
Présentation générale de l'algorithme d'enrichissement de données à l'aide de *Wikidata*

### 3.2.2 Comment traduire des descriptions textuelles datant du XIX<sup>ème</sup> s. en chaînes de caractères pouvant retourner un résultat sur *Wikidata* ?

Dans la réalisation de cet algorithme, la principale difficulté porte sur la récupération d'identifiants *Wikidata* à l'aide de recherches en plein texte. Le script prend en entrée un nom et sa description – tels qu'elles figurent dans des catalogues datant majoritairement du XIX<sup>ème</sup> s.. La difficulté, au delà de la détection et de l'extraction d'informations, est de traduire ces informations pour qu'elles permettent de trouver des résultats pertinents sur *Wikidata*. Ce problème est autant linguistique de technique. Une personne ou une chose est nommée ou décrite d'une certaine manière dans un catalogue de vente ancien. Il n'y a aucune garantie que cette caractérisation corresponde à celle faite par *Wikidata* : l'orthographe des noms évoluent, tout comme la manière de nommer certains métiers. À ces évolutions orthographiques s'ajoutent des évolutions intellectuelles : les titres de noblesse sont un marqueur plus important au XIX<sup>ème</sup> s. français que dans un XXI<sup>ème</sup> s. mondialisé. Une personne n'est que rarement décrite par son titre dans *Wikidata*.

#### Le problème de la traduction des noms

Il existe bien sûr des cas simples, comme l'exemple 3.2 : en extrayant le contenu du `tei:nameet` en traduisant le « roi » issu du `tei:trait`, la chaîne de caractère obtenue est « Henri IV king ». En recherchant cette chaîne de caractère sur *Wikidata*, le premier résultat obtenu est correct. Cependant, de nombreux cas sont plus complexes, surtout lorsque l'auteur.ice du manuscrit est moins célèbre. L'exemple 3.3 est éclairant : dans le catalogue, la personne est nommée « Bruno Daru » ; sur *Wikidata*, le nom de la personne est « Pierre Daru », et son nom complet Pierre Antoine Noël Bruno Daru. Si la recherche en plein texte est faite avec les mêmes paramètres que pour l'exemple précédent (nom de la personne et titre de noblesse), le premier résultat obtenu n'est pas le bon : c'est un renvoi à un article de *l'Encyclopédia Britannica* datant de 1911. C'est en cherchant seulement le nom est le prénom que *Wikidata* retourne un résultat pertinent. Deux choses sont à retenir de cet exemple : dans les catalogues, le prénom d'une personne correspond en fait souvent à son deuxième ou troisième prénom ; ensuite, le titre de noblesse est un critère plus fréquemment mentionné dans les catalogues que dans *Wikidata*. Cela s'explique assez aisément : le XIX<sup>ème</sup> s. connaît une alternance de régimes politiques (royauté, empire, république) où la noblesse n'a pas encore perdu son pouvoir. La probabilité qu'un titre de noblesse soit mentionné sur *Wikidata* diminue lorsqu'un titre est peu important ; dans les catalogues, cependant, même les titres les moins importants sont régulièrement mentionnés. Par conséquent, seuls les titres les plus importants seront extraits pour faire des recherches sur *Wikidata*.

Dans le cas de noms de personnes étrangères, la situation peut être plus complexe



```

1 <item n="134" xml:id="CAT_000233_e134">
2   <!-- ... -->
3   <name type="author">Henri IV</name>
4   <trait>
5     <p>roi de France.</p>
6   </trait>
7   <!-- ... -->
8 </item>

```

Code source 3.2 – Un cas simple : Henri IV roi de France

```

1 <item n="98" xml:id="CAT_000082_e98">
2   <!-- ... -->
3   <name type="author">Daru (Bruno, comte)</name>
4   <trait>
5     <p>célèbre ministre de Napoléon Ier, historien de Venise, de
      ↪ l'Acad. fr., né à Montpellier</p>
6   </trait>
7   <!-- ... -->
8 </item>

```

Code source 3.3 – Un cas plus complexe : Pierre Antoine Noël Bruno Daru

encore. L'exemple 3.4 combine différentes difficultés.

- D'abord, la personne est étrangère ; dans les catalogues, les noms sont systématiquement françaisés – « Albert-Venceslas-Eusèbe » dans le catalogue, « Albrecht Wenzel Eusebius » en langue originelle. Se pose donc la question de si le nom doit être traduit, et si oui comment ?
- Ensuite, comme l'indique la présence de « dit » dans le `tei:name`, il est mentionné un nom de naissance (« de Waldstein ») et un nom d'usage (« Wallenstein »). Idéalement, il faudrait choisir entre l'un ou l'autre, plutôt que de rechercher « Waldstein Wallenstein » sur *Wikidata*, ce qui risque d'augmenter le bruit.

Notre approche s'appuyant sur la structure du texte, le deuxième point peut être réglé : le nom d'usage est écrit au début, et le nom de naissance entre parenthèses (c'est également le cas des noms de personnes nobles, par exemple). Il est donc possible de choisir l'un ou l'autre nom. Le premier point est plus problématique : si la traduction du nom serait envisageable en théorie, celle-ci est difficilement compatible avec une approche basée sur la détection de motifs dans le texte : un prénom est repérable comme étant un motif (trois noms séparés par des tirets) ; cependant, il est impossible de le traduire automatiquement (ce qui demanderait de connaître la langue dans laquelle un prénom doit être traduit). C'est ici que les informations contenues dans le `tei:trait` prennent

```

1 <item n="5518" xml:id="CAT_000401_e5518">
2   <!-- ... -->
3   <name type="author">Wallenstein (Albert-Venceslas-Eusèbe de
   ↪ Waldstein dit)</name>
4   <trait>
5     <p>duc de Friedland, célèbre général de la guerre de Trente ans.
   ↪ Assassiné en 1634.</p>
6   </trait>
7   <!-- ... -->
8 </item>

```

Code source 3.4 – Le problème des noms de personnes étrangères

leur importance : lorsqu'il y a des défaillances dans les informations nominatives, des données biographiques permettent de diminuer le risque d'erreurs. Dans cet exemple, recherche « Albert-Venceslas-Eusèbe Waldstein » ne retourne aucun résultat, de même que rechercher Albert-Venceslas-Eusèbe Wallenstein. Cependant, le bon résultat est obtenu en recherchant « Wallenstein 1634 ». Une difficulté supplémentaire vient avec ce type de cas : différents paramètres de recherche (nom, prénom...) ont un impact différent dans l'obtention du bon résultat en fonction des personnes sur qui la requête est faite. Dans ce cas, rechercher le nom d'usage et la date de naissance retourne un résultat valide, ce qui n'est pas toujours le cas. Pour contourner ce problème, trois solutions ont été mises en place : d'abord, ce types de requêtes a été fait « à la main », de façon non-automatique, pour de nombreuses entrées différentes afin de déterminer la meilleure combinaison de caractères ; ensuite, des tests qui permettent de mesurer l'influence de chaque paramètre de recherche dans l'obtention du résultat ont été réalisés ; enfin, l'algorithme final lance successivement différentes requêtes avec différents paramètres afin de maximiser la probabilité d'obtenir un résultat valide. Nous reviendrons plus en détail sur les deux derniers points.

### L'extraction d'informations biographiques : une autre difficulté

Cependant, le problème ne s'arrête pas qu'aux noms. Dans un exemple ; précédent, le titre de noblesse influençait l'obtention d'un résultat valide. De nombreuses autres informations biographiques pourraient, au premier abord, permettre d'obtenir le bon résultat. C'est souvent le cas, puisque extraire le métier ou la fonction d'une personne permet de supprimer les faux positifs retournés par l'API. C'est par exemple le cas dans l'exemple 3.5. En cherchant uniquement le nom et le prénom (« Hans Bulow »), le premier résultat retourné renvoie à un journaliste suédois. Extraire le mot « pianiste » `tei:trait` et le traduit en anglais permet d'obtenir le bon résultat.

L'extraction d'informations biographiques et leur utilisation dans des requêtes est

```

1 <item n="136" xml:id="CAT_000189_e136">
2   <!-- .. -->
3   <name type="author">Bulow (Hans)</name>
4   <trait>
5     <p>le célèbre pianiste.</p>
6   </trait>
7 </item>

```

Code source 3.5 – Un exemple où l’extraction du métier permet l’obtention du bon résultat

donc pertinent. Cependant, extraire trop d’informations conduit à lancer des requêtes qui ne renvoient aucun résultat. Dans les exemples 3.6 et 3.7, extraire et traduire des fonctions conduit à lancer les requêtes « John Okey colonel » et « Jean Bouhier président » qui ne retournent aucun résultat, ou des résultats qui ne sont pas valides. Cependant, dans les deux cas, si une requête est lancée sans le métier, un résultat correct est obtenu. Les raisons pour lesquelles des résultats erronés sont retournés ne sont cependant pas les mêmes, et il est intéressant de mieux observer les requêtes lancées et les résultats obtenus. Dans le premier cas, le terme mis en avant dans le `tei:trait` (« colonel ») n’est pas celui avec lequel la personne est décrite sur *Wikidata* (où John Okey est décrit comme étant un homme politique). Cela met en avant un problème relatif au changement de regard sur des personnalités : dans un contexte, la personne est décrite comme une figure militaire, dans l’autre comme une figure politique. Le deuxième cas est plus technique. Il y a en fait une erreur dans la requête qui ne retourne pas de résultat (« Jean Bouhier président ») : un.e président.e de parlement n’est en général pas décrite comme « président ». Cependant, en extrayant des données uniquement par détection de motifs, il est possible de repérer et traduire un terme générique comme « président ». Extraire le complément « Parlement de Dijon » du `tei:trait` n’est cependant pas possible (cela impliquerait d’étudier la grammaire de la phrase, pour mettre en avant la relation entre « président » et « Parlement de Dijon »). Au vu de la taille et de la variété du jeu de données, il est impossible de traiter au cas par cas des entrées, ou préciser suffisamment la détection de motifs pour pouvoir résoudre ce genre de difficultés.

De situations comme les exemples 3.6 et 3.7, il faut donc retenir que l’extraction d’informations vient nécessairement avec un risque d’erreur. Le parti pris a donc été de ne pas repérer les métiers et autres termes très spécifiques, comme les grades militaires et les titres de noblesse peu élevés : ils ne retournent pas de résultats sur le moteur de recherche. Ensuite, plus des requêtes sont précises, plus elles risquent de retourner du silence (c’est-à-dire, de ne pas donner de réponse) ; cependant, si un résultat est obtenu, il est plus probable que ce résultat soit correct. Une fois l’extraction d’informations faite, l’algorithme de résolution d’entités avec *Wikidata* a donc été conçu en suivant un principe soustractif : les premières recherches sont faites avec un maximum de paramètres ; si

```

1 <item n="152" xml:id="CAT_000189_e152">
2   <!-- ... -->
3   <name type="author">Okey (John)</name>
4   <trait>
5     <p>colonel anglais, un des lieutenants de Cromwell.</p>
6   </trait>
7   <!-- ... -->
8 </item>

```

Code source 3.6 – Quand l’extraction d’un métier conduit à des requêtes trop spécifiques

```

1 <item n="5430" xml:id="CAT_000401_e5430">
2   <!-- ... -->
3   <name type="author">Bouhier (Jean)</name>
4   <trait>
5     <p>président au Parlement de Dijon, membre de l'Académie
6       ↪ française.</p>
7   </trait>
8   <!-- ... -->
9 </item>

```

Code source 3.7 – Le cas des métiers dont l’extraction est problématique

aucun résultat n’est obtenu, des paramètres sont enlevés pour que l’API retourne un plus grand nombre de résultats. Enfin, ces deux exemples montrent qu’il n’est pas possible d’extraire et de traduire des informations sans prendre en compte ce qui sera pertinent pour le moteur de recherche de *Wikidata*. Il ne s’agit donc pas seulement d’extraire des informations, mais aussi de s’adapter avec ce moteur de recherche pour augmenter la probabilité d’obtenir un résultat valide.

### 3.2.3 Comment négocier avec le moteur de recherche de *Wikidata* ?

Comme cela commence à apparaître, l’extraction d’informations, lorsqu’elle vise à interagir avec des données externes, vient avec des difficultés supplémentaires. Il ne faut pas seulement extraire les informations ; leur extraction et structuration doivent permettre de lancer des recherches en plein texte, et donc de minimiser le bruit (les informations non pertinentes) et le silence (l’absence d’informations) de la part du moteur de recherche. Il faut donc traduire les informations extraites pour qu’elles correspondent au vocabulaire utilisé par *Wikidata*. Cette opération n’est pas anodine : si les catalogues de vente fonctionnent avec leurs propres catégories, le même peut être dit de *Wikidata* : certains types de données sont plus souvent référencées que d’autres et *Wikidata* utilise

un vocabulaire qui lui est propre. Pour bien mener ce processus de traduction et de structuration de l'information, il est nécessaire de bien connaître le fonctionnement de ce moteur de recherche pour mieux s'y adapter.

Comme cela a été dit, l'alignement avec *Wikidata* passe par l'utilisation de l'API mise en point par l'institution afin de lancer automatiquement des recherches en plein texte ; l'objectif est que le premier résultat retourné par le moteur de recherche soit le bon. La première chose à remarquer est que, contrairement à un moteur de recherche généraliste (comme *Google*, *QWant...*), ce moteur n'est pas compatible avec des requêtes approximatives. L'exemple 3.8 est pertinent à ce égard<sup>19</sup>. Dans de nombreuses entrées, comme c'est le cas ici, les fonctions d'une personne ayant participé à la révolution sont présentées de façon précise : Marc David Alba Lasource est décrit comme étant un « conventionnel girondin ». Cette mention, régulièrement présente dans les catalogues, pourrait être relevée en tant que telle. Cependant, lancer la recherche « Lasource conventionnel » ne retourne aucun résultat. Si la même recherche est lancée sur un moteur de recherche généraliste (ici, *QWant*), la page *Wikipedia* de Lasource fait partie des premiers résultats<sup>20</sup>. Cette différence dans les données retournées par les moteurs de recherche a deux explications : un moteur de recherche généraliste recherche les occurrences de mots, non seulement dans le titre de la page, mais aussi dans le corps du texte. Si le mot « conventionnel » est absent du titre, il est certainement à plusieurs reprises dans une notice biographique type *Wikipedia*. *Wikidata* ne contenant que des données, et pas de texte en tant que tel, l'indexation du corps du texte par le moteur de recherche interne à *Wikidata* n'est pas possible. Ensuite, la plupart des moteurs de recherche généralistes utilisent des méthodes de traitement du langage afin de simplifier la requête lancée par l'utilisateur.ice : les mots recherchés sont simplifiés, le moteur de recherche associe les termes recherchés avec d'autres termes « cooccurrents », c'est-à-dire fréquemment utilisés ensemble<sup>21</sup>. Dans le cas du moteur de recherche de *Wikidata*, la requête de l'utilisateur.ice ne semble pas être retraitée : des signes de ponctuation ou des fautes de frappe influencent l'obtention d'un résultat, de même que l'usage de termes inadaptés.

Pour faire face à la « rigidité » relative du moteur de recherche de *Wikidata*, il est

---

19. Dans cet exemple, le prénom, « M.-D.-A. », n'est pas pris en compte pour se concentrer sur l'utilisation d'informations biographiques dans le `tei:trait`.

20. Le 29/07/2022, c'est le troisième résultat ; le premier correspond à une vente aux enchères d'archives du conventionnel. Les moteurs de recherche pouvant être mis à jour régulièrement, il est possible que l'ordre des résultats change

21. *Moteur de recherche*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Moteur\\_de\\_recherche](https://fr.wikipedia.org/wiki/Moteur_de_recherche) (visité le 17/07/2022). Pour des analyses plus détaillées sur la construction d'ensemble de termes cooccurrents via le développement de vecteurs de mots, voir Tomas Mikolov, Kai Chen, Greg Corrado et Jeffrey Dean, « Efficient Estimation of Word Representations in Vector Space » (, 2013), Publisher : arXiv Version Number : 3, DOI : 10.48550/ARXIV.1301.3781 ; pour un article technique détaillant la classification et la sélection de résultats pertinents par apprentissage profond, voir Paul Covington, Jay Adams et Ernie Sargin, « Deep Neural Networks for YouTube Recommendations », dans *Proceedings of the 10th ACM Conference on Recommender Systems*, New York, 2016, URL : <https://research.google/pubs/pub45530/> (visité le 10/06/2022)

```

1 <item n="140" xml:id="CAT_000197_e140">
2   <!-- ... -->
3   <name type="author">Lasource (M.-D.-A.)</name>
4   <trait>
5     <p>célèbre conventionnel girondin, né près de Montpellier en
6       ↪ 1762, guillotiné en 1793.</p>
7   </trait>
8   <!-- ... -->
9 </item>

```

Code source 3.8 – Le problème de l’approximation et de la traduction : Lasource, conventionnel

donc nécessaire de préparer ses données au moment de leur extraction. En prenant le même exemple (3.8), un résultat correct peut être obtenu en remplaçant « conventionnel » par « politician »<sup>22</sup>, pour rechercher sur *Wikidata* « Lasource politician ». Ici, la traduction de « conventionnel » en « politician » est d’autant plus intéressante que la date de naissance dans le catalogue (1762) ne correspond pas à celle indiquée sur *Wikidata* (1763). Dans un cas comme celui-ci, où certaines données sont incorrectes, il est important d’extraire un maximum d’informations pour que, si certaines requêtes ne rapportent pas de résultats, pouvoir en faire d’autres avec différents paramètres.

En conclusion, il faut retenir que le moteur de recherche de *Wikidata* n’admet pas d’erreurs, ni de requêtes partiellement erronées (dans l’exemple 3.8, où la date de naissance soit correcte, mais pas la date de décès) ; il ne prend pas non plus en compte la synonymie, ce qui veut dire qu’il n’améliore pas la requête lancée par un.e utilisateur.ice... Cela signifie que les termes utilisés dans une requête doivent être adaptés à ceux que *Wikidata* utilise. Les termes spécifiques utilisés dans les catalogues (« conventionnel »), mais aussi de nombreux titres militaires et de noblesse peu élevés (« capitaine », « marquis ») sont relativement rarement présents sur *Wikidata*. Lorsque les requêtes sont lancées, de tels termes sont donc abandonnés et parfois remplacés par des termes plus génériques : par exemple, « capitaine » est remplacé par « military », traduction anglaise de « militaire ». De même, des termes principalement en usage dans la langue française, comme « conventionnel » sont moins efficaces pour lancer des recherches.

### 3.3 Une approche prédictive

L’alignement avec *Wikidata* et l’extraction d’entités n’est donc pas une opération anodine : les données contenues dans les catalogues sont variées, autant par leur structure que par les informations qu’elles contiennent ; il peut être difficile à faire la traduction

---

22. « Personnalité politique »

1

```
<name type="author">Verneuil (Charlotte Séguier duchesse  
↪ de)</name>
```

Code source 3.9 – Peut-on identifier les différents éléments d’une phrase par détection de motifs ?

de données du XIX<sup>ème</sup> s. en chaînes de caractères pouvant retourner des réponses valides sur *Wikidata* ; enfin, le l’alignement repose sur une bonne connaissance du moteur de recherche de *Wikidata*.

De plus, la technique utilisée dans l’extraction de données, reposant sur la détection de motifs à l’aide de expressions régulières et de tables de conversion, est une technique qui vient avec un certain nombre d’incertitudes. Avec ce genre de techniques, il est impossible de « comprendre » ce qu’un élément signifie. Dans l’exemple 3.9, formellement, rien ne sépare le nom propre de la duchesse (« Séguier ») du nom de son duché (« Verneuil »). En s’appuyant sur une connaissance de la structure répétitive des entrées, il est uniquement possible de supposer que le nom entre parenthèses est un nom propre, tandis que le nom hors des parenthèses correspond au nom du duché. En bref, les méthodes de détection de motifs utilisées, peuvent uniquement inférer le sens d’un mot par rapport à sa position dans une phrase. Si cette technique implique une certaine incertitude, elle est cependant particulièrement adaptée à un corpus semi-structuré, comme c’est le cas des catalogues de vente de manuscrits, et à l’opération d’alignement avec *Wikidata*. Comme cela est expliqué plus bas, au fond, il n’est pas tellement important de distinguer le sens des différentes informations : ce qui a du sens, c’est que l’extraction et la structuration des informations permet de construire des chaînes de caractères à rechercher sur *Wikidata*. Identifier la fonction d’un mot n’est donc ici qu’un moyen – contrairement à de l’analyse lexicale, où la fonction des mots dans une phrase est signifiante –. En effet, repérer le rôle que tiennent les termes extraits (métier, prénom...) permet de mieux construire la chaîne de caractère recherchée sur *Wikidata*, en pouvant filtrer certaines informations (retirer les dates de vie et de mort, par exemple).

Étant donnée cette quantité d’incertitudes, l’approche suivie dans l’alignement avec *Wikidata* peut être qualifiée de « prédictive ». Par ce terme, il faut comprendre que il n’y a pas, de certitude totale dans le processus d’extraction et de traduction des données. Il n’est pas possible de récupérer avec une certitude totale le bon identifiant. L’objectif cet algorithme n’est donc pas de trouver la « bonne » réponse. Il est de construire une chaîne de caractère dont on prédit qu’elle apportera un résultat pertinent. De la même manière, la phase de préparation des données est un processus qui sélectionne et normalise certaines informations dont on considère – après un long processus de test et d’essais – qu’elles seront pertinentes dans l’obtention des bons résultats. Enfin, le premier rôle des tests est de quantifier les prédictions. Ils répondent à la question : étant donné les

résultats obtenus lors des tests, quelle est la probabilité que la prochaine chaîne de caractères recherchée retourne un résultat pertinent ? Cette approche prédictive implique nécessairement un degré d'incertitude, et donc le développement d'algorithmes flexibles qui cherchent à minimiser le bruit.

Être conscient de la nature prédictive de ce processus et quantifier la qualité des algorithmes à l'aide de tests permet cependant de prendre de meilleures décisions techniques. La lecture distante et la détection de motifs supposent d'avancer « à l'aveugle », en s'appuyant sur sa connaissance de la structure du texte pour extraire les bonnes informations. Étant donné qu'il est impossible d'être totalement certain que les bonnes données ont été extraites, l'étape suivante – le lancement des requêtes sur l'API – doit malléable et s'adapter aux données disponibles. C'est pourquoi le parti pris a été de concevoir un algorithme qui continue de lancer des requêtes en retirant des paramètres tant qu'un identifiant n'a pas été trouvé.



# Chapitre 4

## Un algorithme de détection de motifs pour préparer et structurer les données

Avant de chercher à récupérer un identifiant *Wikidata* via l'API, un algorithme se charge de traduire et de structurer les données : à partir d'un nom et de son éventuelle description, un dictionnaire qui contient les informations de manière structurée est construit. Cette étape était initialement censée être une simple extraction d'information : à partir du `tei:name` et du `tei:trait`, un ensemble d'informations étaient mises bout à bout afin de former une chaîne de caractères à rechercher sur l'API. Le processus s'est complexifié pour intégrer l'extraction, la traduction et la structuration des données. En construisant un dictionnaire à partir de texte, il est possible de savoir précisément quelles données sont disponibles pour lancer des requêtes ; plusieurs requêtes peuvent alors être lancées sur l'API avec différents paramètres, ce qui permet d'augmenter les probabilités d'obtenir un identifiant valide.

### 4.1 Présentation générale

#### 4.1.1 Les formats d'entrée et de sortie

Le but de l'extraction de données permet de transformer la représentation TEI visible en 4.1 – représentée sous forme d'un TSV pour faciliter la lecture des données – au dictionnaire visible en 4.2. Voici la signification des différentes clés<sup>1</sup> du format de sortie :

- **fname** : cette clé permet d'accéder au prénom d'une personne. Les données contenues dans cette clé viennent du `tei:name`. **fname** est l'abréviation de « first name ».

---

1. Une clé de dictionnaire est l'élément à gauche des « : » ; la clé permet d'accéder à la valeur, visible à droite du « : », ce qui permet d'associer des valeurs entre elles, et donc de stocker des objets ou de remplacer une clé présente dans un texte par une valeur, par exemple.

```

1 <item n="271" xml:id="CAT_000327_e271">
2   <!-- ... -->
3   <name type="author">Turenne (Henri de La Tour d'Auvergne vicomte
4     ↪ de)</name>
5   <trait>
6     <p>illustre maréchal de France, né en 1611, tué en 1675.</p>
7   </trait>
8   <!-- ... -->
9 </item>

```

Code source 4.1 – L’entrée XML-TEI à partir de laquelle des données sont extraites

- **lname** : cette clé permet d’accéder au nom de famille de quelqu’un. C’est cette information, extraite du `tei:name`, qui est centrale aux requêtes. **lname** abréviation de « last name »)
- **nobname\_sts** : cette clé contient un nom de famille noble. Dans ces cas, un titre de noblesse est présent dans les entrées de catalogue (seuls les titres de noblesse les plus importants sont extraits du dictionnaire, ce qui n’est pas le cas ici). Les informations contenues ici proviennent du `tei:name`. Cette clé est l’abréviation de « nobility name\_status »
- **status** : le statut d’une personne, soit son titre de noblesse (ici, le titre « vicomte » n’a pas été extrait car il est rarement présent sur *Wikidata*). Les informations contenues ici proviennent en général du `tei:name` et parfois du `tei:trait`.
- **dates** : les dates de naissance ou de mort d’une personne (seules ces dates sont conservées). Ces informations proviennent du `tei:trait`.
- **function** : la fonction d’une personne, soit, en général, son métier ou son occupation principale. Cette information provient du `tei:trait`.
- **rebuilt** : un booléen indiquant si un prénom a été reconstruit à partir d’initiales ou non.

Comme cela a été dit auparavant, le nom attribué à ces clés n’est pas systématiquement indicateur des valeurs qui y sont associées : si l’entrée de catalogue correspond à une personne, alors les clés correspondent aux informations qu’elles contiennent. Si l’entrée de catalogue ne correspond pas à une personne, ces clés seront également utilisées. Ce qui est important, c’est la hiérarchie d’importance entre les différentes clés : **lname** est la clé centrale et contient presque toujours des informations, **fname** des données secondaires et **date** des dates. Les autres clés sont rarement utilisées si l’entrée de catalogue ne correspond pas à une personne.

```
1 {  
2   "fname": "henri ",  
3   "lname": "la tour d'auvergne",  
4   "nobname_sts": "Turenne ",  
5   "status": "",  
6   "dates": "1611 1675 ",  
7   "function": "marshal",  
8   "rebuilt": False  
9 }
```

Code source 4.2 – La sortie JSON correspondante

### 4.1.2 Présentation de l’algorithme d’extraction d’informations

L’algorithme détaillé ci dessous est présenté sous forme graphique dans la figure 4.1. Cette étape peut être séparée en deux parties différentes : l’extraction d’informations nominatives du `tei:name` et la récupération de données biographiques du `tei:trait`.

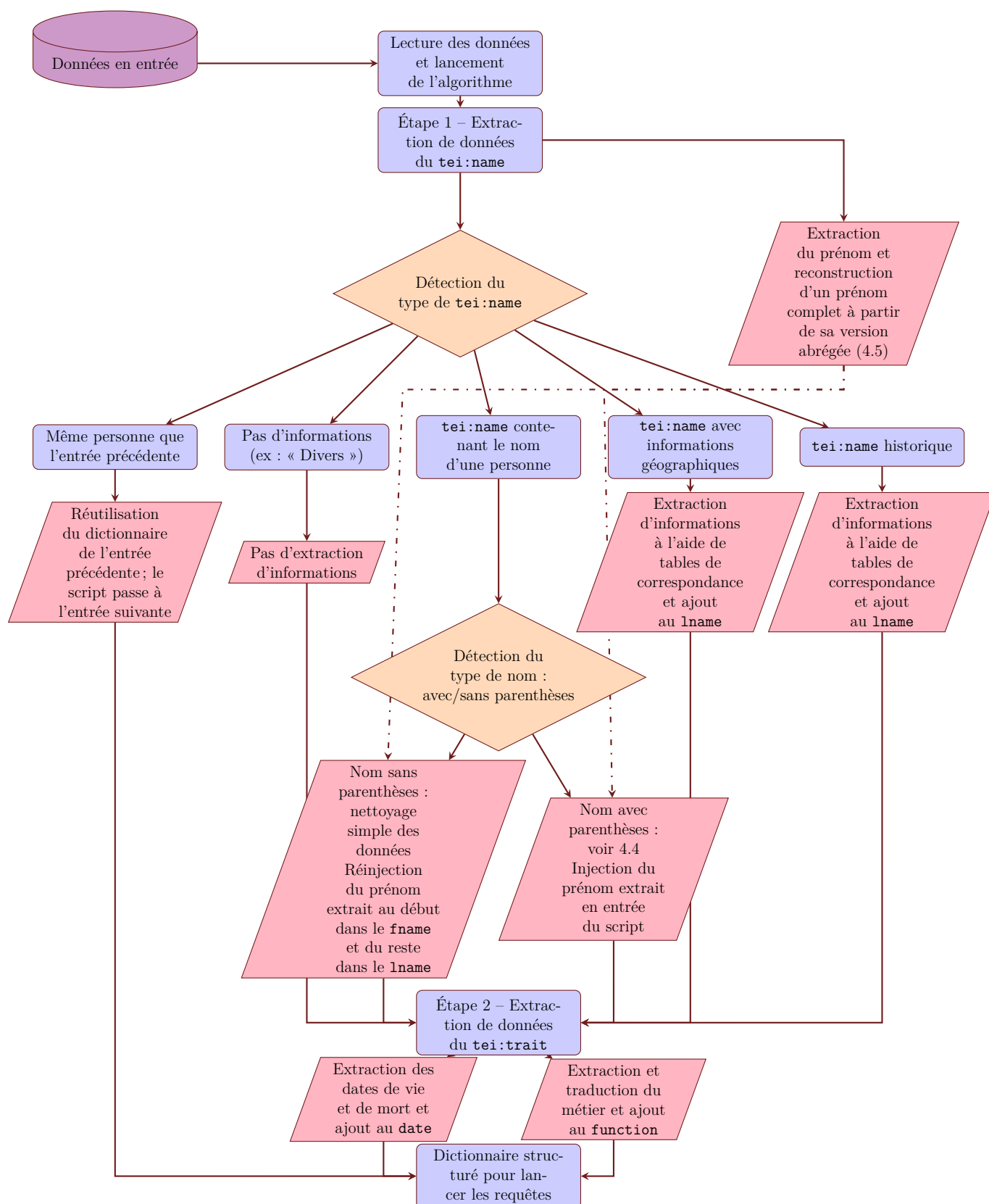
L’extraction d’informations du `tei:name` est l’étape plus complexe. La difficulté tient au fait que cette balise peut contenir des informations variées et structurées de façon très différente. Cela demande d’identifier des motifs récurrents et de les repérer dans le texte à différents degrés et à différentes étapes. Dans un premier temps, les prénoms sont systématiquement extraits. Ils sont toujours détectés – même si la balise ne contient le nom d’une personne : cette extraction permet justement d’identifier le type de données contenues dans le `tei:name`. Les prénoms sont repérés à l’aide de plusieurs expressions régulières qui permettent d’identifier des prénoms complets et abrégés, qu’ils soient composés ou non. Cette détection prend en compte les différents types d’abréviations possibles (un prénom composé peut être entièrement abrégé ; à l’inverse, seulement un des prénoms peut être abrégé) et les différentes typographies (séparer les prénoms avec des traits d’union ou non, par exemple). Dans le cas où un prénom serait abrégé, il est si possible reconstruit : des initiales sont remplacées par un nom complet ; ce processus est présenté plus en détail ci-dessous. Cela permet d’augmenter le taux de réussite dans l’alignement avec des identifiants *Wikidata*, mais vient avec plusieurs difficultés techniques, comme nous le verrons. Une fois ce nom extrait, le type d’information contenue dans le `tei:name` doit être identifié : en fonction du type d’information (géographique, historique, nominative...), différents traitements sont mis en place. Cette identification se fait par détection de motifs augmentée par l’usage de tables de conversion<sup>2</sup> et de listes contenant du vocabulaire spécifique. Listes et tables étant classées thématiquement, il est possible, par un processus éliminatoire, d’identifier avec certitude le type d’information contenue dans le `tei:name`. Le traitement du `tei:name` dépend grandement de cette détection : si cette balise contient des éléments géographiques ou historiques, l’extraction d’informa-

---

2. Pour un exemple de table de conversion, voir 1

tions repose en grande partie sur les tables de conversion. S'il s'agit en revanche d'un nom de personne, il est alors nécessaire d'identifier les différents types de données nominatives (prénom, nom de famille, nom de famille noble...) pour bien structurer les données. En effet, c'est de cette structure que dépend la bonne construction du dictionnaire, et donc la constitution de requêtes adaptées à l'API. Ce processus d'extraction des données s'appuie majoritairement sur une détection de motifs. Le motif déterminant est la présence ou non dans parenthèses dans le nom. Comme nous le verrons, si un nom contient des parenthèses, les informations sont bien plus structurées que s'il n'en contient pas. Les informations peuvent alors être extraites avec une bien plus grande granularité.

Une fois les informations nominatives extraites du `tei:name`, il reste à extraire les données biographiques pertinentes du `tei:trait`. Cette étape, plus simple que la précédente, vaut principalement pour les entrées où c'est l'auteur.ice d'un document qui est mentionné.e dans le `tei:name`. Les seules informations extraites concernent les dates de vie et de mort des personnes, ainsi que son métier. Quelques difficultés techniques subsistent cependant : il faut notamment distinguer une date de naissance/décès d'une autre date, afin de diminuer le bruit ; de plus, il faut réussir à extraire de façon automatique l'occupation principale d'une personne, lorsque plusieurs personnes sont mentionnées (« militaire » et « auteur », par exemple). La résolution de ces deux problèmes repose sur une bonne connaissance du corpus de catalogues et de la structure des `tei:trait`.

FIGURE 4.1 – Processus d'extraction d'informations du `tei:name` et `tei:trait`

## 4.2 Identifier le type de nom

Les éléments `tei:name` contiennent le titre donné à l'item vendu. Si c'est souvent le nom de l'auteur.ice du document, ce n'est pas toujours le cas. L'extraction d'éléments du `tei:name` dépend, comme cela a été dit, de l'identification du « type » de nom. La décision a été prise de classer tous les `tei:name` en cinq catégories.

Les noms génériques : ces éléments ne contiennent pas d'informations précises. Les entités *Wikidata* étant spécifiques plutôt que génériques, il n'est pas certain que les entrées puissent être alignées avec *Wikidata*; si des entités *Wikidata* correspondent à ces éléments, les informations qu'elles contiennent seront probablement trop génériques pour être utilisables dans un contexte économétrique. Lorsque cela est possible, l'alignement est quand même fait; c'est le cas par exemple pour les `tei:name` contenant la mention de chartes. Dans ces cas, un dictionnaire vide est retourné et l'alignement avec *Wikidata* n'aura pas lieu. Une exception est faite dans le cas chartes, où un alignement est fait avec l'entité *Wikidata* génériques. Dans cette catégorie se trouve par exemple :

— `<name type="author">DIVERS</name>`

Les noms de type « Le même » ou « La même »; lorsqu'il y a cette valeur dans le `tei:name`, l'auteur.ice est la même personne que l'auteur.ice de l'entrée de catalogue précédente. Dans ce cas, le dictionnaire de cette entrée est réutilisé. Par exemple :

1. `<name type="author">Le même</name>`.

Les noms géographiques; ces entrées sont détectées à l'aide de tables de conversion et de listes. Celles-ci sont classées thématiquement : liste d'anciennes colonies françaises (4), de départements français du XIX<sup>ème</sup> s.(2), d'anciennes provinces françaises (5) et de pays (3). Dans ce cas, un alignement avec une entité *Wikidata* est possible; cependant, il n'est pas toujours envisageable de retrouver l'entité précise. En effet, un `tei:name` peut contenir une mention d'une donnée géographique, mais également d'autres détails. C'est par exemple le cas dans le quatrième exemple ci-dessous. Il faut alors aligner le `tei:name` avec son équivalent générique sur *Wikidata* (l'exemple ci-dessous, par exemple, a été aligné uniquement avec l'entité « Paris »). Dans cette catégorie se trouvent :

— `<name type="other">AISNE (département de 1')</name>`

— `<name type="author">Bourbonnais. </name>`

— `<name type="author">Paris : Musée royal du Louvre</name>`

— `<name type="author">Garde nationale parisienne en 1792 (brevet de volontaire de 1a)</name>`

Les noms correspondant à des événements historiques. Là encore, une table de conversion est utilisée (6). Ici, une difficulté apparaît cependant : du fait de la variété des événements historiques mentionnés dans les entrées de catalogue, il n'est pas possible

d'enregistrer l'ensemble des évènements dans des tables afin de permettre une détection de tous les évènements. Les `tei:name` ont donc été analysés pour extraire les évènements les plus importants. Ensuite, comme pour les termes géographiques, il n'est pas possible de donner aux tables de conversion une granularité suffisamment fine pour contenir toutes les données possibles. Des alignements partiels ont donc été faits : le premier exemple ci-dessous a été aligné avec l'entité *Wikidata* « Révolution française ».

- `<name xmlns="http://www.tei-c.org/ns/1.0" type="other">THÉÂTRE  
RÉVOLUTIONNAIRE</name>`
- `<name type="author">COMMUNE DE 1871.</name>`
- `<name type="other">Siège de La Rochelle en 1628</name>`

Les noms de personnes. Ceux-ci ne sont pas simples à traiter : ils peuvent contenir de nombreuses informations : deux noms de famille (usuels et nobles), titres de noblesse, plusieurs prénoms. Ils ont également une structure variée, comme cela apparaît dans les exemples ci-dessous : les noms peuvent être écrits en utilisant des parenthèses ou non ; un prénom peut être écrit en entier, comme dans premier exemple, entièrement abrégé (comme dans l'exemple 3) ou encore partiellement abrégé, ce qui est le cas dans le troisième exemple. Ces différences, qui ne posent pas de problème à un regard humain, sont autant de problèmes techniques. En effet, la détection de motifs fonctionne uniquement sur des critères formels, ou structurels. Il faut donc, en s'appuyant sur la structure des documents, réussir à distinguer un prénom d'un nom de famille, un nom de famille d'un autre, ou encore un nom abrégé d'un nom complet.

- `<name type="author">Humboldt (le baron Alexandre de)</name>`
- `<name type="author">Taccani Tasca (madame la comtesse)</name>`
- `<name type="author">LEGOUVÉ (G. M. J. B.)</name>`
- `<name type="author">LOUIS XVIII</name>`
- `<name type="author">Duras (Emm.-F. de Durfort, duc de)</name>`

Si l'on pose l'extraction d'informations nominatives sur une personne comme étant l'objectif principal, une difficulté apparaît vite : qu'est-ce qui distingue un nom de personne d'un des autres types de noms ? Les éléments `tei:name` sont voués à contenir des noms propres ; chercher à distinguer les noms propres des noms communs n'a donc pas d'intérêt. Cela est d'autant plus que la graphie varie d'un catalogue à l'autre : dans certains, les majuscules sont signifiantes et pourraient permettre de distinguer noms communs et noms propres, tandis que dans d'autres, l'intégralité du `tei:name` est en majuscule. Une détection de motifs à partir de simples critères formels (du type : « Un nom de personne est un ou plusieurs mots commençant par des majuscules ») n'est pas opérante pour ce corpus. Il n'est pas non plus possible de définir un nom positivement, puisqu'il n'existe aucun

critère définitoire pour un nom propre. Enfin, à cette étape, il n'est pas non plus possible de s'appuyer sur l'extraction de prénoms. L'extraction de prénoms se base, comme on le verra, sur de la détection de motifs ; là encore, ce qui est identifié comme un prénom peut tout aussi bien être un nom propre, et il n'existe à ce stade aucune possibilité pour distinguer un nom propre d'un autre. C'est à ce stade qu'a été décidée l'utilisation de tables de conversion et de listes contenant du vocabulaire thématique : en identifiant des références récurrentes à des événements ou des lieux dans les `tei:name`, il devient possible de définir les noms de personne négativement. En définitive, un nom de personne, c'est donc ce qui n'est pas autre chose et la détection du type de nom fonctionne donc de manière éliminatoire.

L'algorithme mène donc une série de tests, cherchant à détecter si un `tei:name` rentre dans telle ou telle catégorie. Du fait du fonctionnement technique de `python`, une série de tests éliminatoires doit aller du cas le plus particulier au cas le plus générique afin d'éviter les faux positifs : une fois qu'un élément a été détecté comme appartenant à une catégorie, il ne peut plus être réassigné à une autre. C'est pourquoi l'algorithme commence par chercher à classer les éléments dans les catégories où le taux d'erreur est le plus faible, pour ensuite finir par la catégorie la plus générique : celle des noms de personne. : c'est dans ces catégories que le taux d'erreur est le plus faible.

Le script commence donc par chercher à classer un `tei:name` dans les catégories où les informations sont écrites plus ou moins toujours de la même manière. Il cherche d'abord à identifier si le `tei:name` correspond à « Le même » ou « La même ». Dans ce cas, le `tei:name` est le même que celui de l'entrée précédente et c'est ce dictionnaire qui est réutilisé. Ensuite, les entrées génériques sont détectées. Comme elles contiennent toujours des informations écrites de la même manière (« Documents divers »), le taux d'erreur est là encore très faible. Ces entrées génériques ne contiennent pas d'informations spécifiques ; si un `tei:name` appartient à cette catégorie, alors l'algorithme n'extrait pas d'informations.

Ensuite, si un `tei:name` n'appartient ni à l'une ni à l'autre catégorie, alors l'algorithme cherche à classer un nom en différentes catégories à l'aide de tables de comparaison et de listes contenant du vocabulaire spécifiques. Ces tables et listes sont classées en deux catégories (données historiques et géographiques) afin de définir des traitements spécifiques ; l'algorithme cherche d'abord à identifier des entrées géographiques, puis historiques. En effet, un bien plus grand nombre de tables contenant des données géographiques existe<sup>3</sup>, ce qui augmente les possibilités d'un classement correct. Enfin, le script cherche à identifier des informations historiques dans un `tei:name` (6). Si une donnée géographique ou historique est repérée, alors équivalents normalisés de cette donnée sont ajoutés au dictionnaire grâce à l'usage de tables de conversions.

Par processus d'élimination, si aucune de ces informations n'a été détectée, alors il

---

3. Nom d'anciennes provinces françaises (5), de départements du XIX<sup>ème</sup> s. (2), d'anciennes colonies (4) et de pays (3)



n'est plus possible de classer le `tei:name` dans aucune autre catégorie. Il est alors considéré que le contenu du `tei:name` est le nom d'une personne. L'extraction d'informations se fait ici plus complexe<sup>4</sup>, comme nous le verrons : le script traite le nom différemment en fonction de sa structure (présence ou non de parenthèses dans le nom), puis extrait d'éventuels titres de noblesse, noms de famille noble et nom de famille usuel. Enfin, il extrait des prénoms et cherche à reconstruire un prénom complet à partir de son abréviation. Le processus étant éliminatoire, il n'y a bien sûr aucune certitude totale que le contenu du `tei:name` soit bel et bien le nom d'une personne ; il n'est cependant plus possible de mieux catégoriser les éléments. Cela n'est pas non plus extrêmement important, puisque cet algorithme de classification a un rôle fonctionnel et n'impacte pas la manière dont les identifiants sont récupérés depuis l'API *Wikidata*. Il permet principalement d'adopter un fonctionnement modulaire en cherchant à détecter des motifs spécifiques dans les `tei:name` en fonction de la catégorie à laquelle ils appartiennent. La détection de motifs étant « aveugle », le traitement qui est fait pour cette catégorie peut être mis à profit de différents types de données : les mêmes motifs peuvent être recherchés et extraits de différents types d'entrées.

À l'issue de cette phase de classification, il est possible de mieux connaître les types de `tei:name` et leur répartition dans le corpus. Les noms de personnes sont très majoritaires dans les `tei:name` : ils représentent 80634 entrées sur un total de 82913, soit 97,25% du corpus. C'est donc pour ce type de données que l'extraction des données a été pensée, ce qui permet d'extraire des informations avec une granularité plus fine qu'avec le reste du corpus. Viennent ensuite les noms de lieux (1406 entrées) et les éléments divers (550). Pour finir se trouvent les événements historiques (232 entrées) et les éléments vides, soit 92 entrées<sup>5</sup>. Bien que, comme cela a été dit, la classification des `tei:name` en différents types soit purement fonctionnelle, ces chiffres permettent d'avoir une meilleure connaissance du corpus, et mieux comprendre à partir de quel type de données se fait l'alignement *Wikidata*. Cette classification montre également d'intérêt de la détection de motifs augmentée de tables de correspondances : bien que ce soit une méthode « légère », elle peut efficacement, sur un corpus semi-structuré, être à la base d'une extraction précise de données.

### 4.3 Le traitement des noms de personnes

Comme cela a déjà été dit, le processus d'extraction d'informations du `tei:name` est relativement simple pour les entrées qui ne contiennent pas de noms de personnes : il s'agit principalement de remplacer des informations non-normalisées par leur équivalent

---

4. Pour une représentation graphique de cette étape, voir 4.4

5. Ces chiffres proviennent d'un calcul réalisé à partir du script développé pour identifier le type d'entrée. Ils ont été produits afin de réaliser une représentation graphique du type d'entrée, visible en annexes (2)

normalisé, à l'aide d'un dictionnaire. Pour les noms de personne cependant, le processus est plus complexe : il demande d'identifier titres de noblesse et différents types de noms ; de plus, de nombreux prénoms sont abrégés pour gagner de la place dans les catalogues. Pour augmenter le taux de réussite de l'alignement avec *Wikidata*, une méthode a été développée pour reconstruire un nom complet à partir de son abréviation. Après avoir présenté les différentes étapes de l'extraction d'informations nominatives, cette partie détaillera ce processus de reconstruction du prénom. Dans ces deux étapes, la nature semi-structurée des documents est centrale, puisqu'elle permet d'identifier avec une grande précision les différentes informations.

### 4.3.1 Trouver des solutions adaptées à différents types de noms

Pour extraire des informations d'un `tei:name` et dans leur un statut ou d'une signification spécifique, il est nécessaire d'identifier ce qui distingue un prénom d'un nom de famille ou d'un titre de noblesse. Une première difficulté apparaît vite : formellement, il n'y a pas nécessairement de différence entre ces informations. Dans l'exemple 4.1, le `tei:name` correspond à :

```
1 <name type="author">Turenne (Henri de la Tour d'Auvergne vicomte
   ↪ de)</name>
```

Code source 4.3 – Le `tei:name` de l'exemple 4.1

Les informations à extraire de cet élément sont visibles dans la figure 4.2. Si chaque mot est pris individuellement, il est difficile de les distinguer l'un de l'autre. Le nom de famille noble, le prénom et le nom de famille usuels débutent tous par une majuscule. Il est ici impossible de trouver un motif distinctif pour un nom de famille ou pour un prénom. Cela est d'autant plus vrai qu'un nom peut être composé ; ici, le nom de famille usuel est composé de deux mots débutant par une majuscule et séparés par un « d' ». Étant donné la taille du corpus, il n'est pas non plus possible d'utiliser des tables de conversions pour identifier ce qui est un nom ou un prénom, comme cela a été fait pour les lieux et les événements géographiques. La seule information qui peut être traitée à l'aide d'une table de conversion est le titre de noblesse, puisque ces titres existent en nombre limité. Par conséquent, aucun élément à l'intérieur des termes à relever ne permet de les identifier. Là où ils peuvent cependant être distingués, c'est au niveau de leur place relative au sein de la phrase. Pour reprendre l'exemple 4.2, le nom de famille noble se retrouve hors de la parenthèse, au tout début de la phrase ; les autres informations sont dans la parenthèse, où le prénom est suivi du nom de famille et du titre. S'il n'est pas possible d'identifier la valeur d'un nom à partir de ses caractéristiques propres, il est donc possible de la déduire de sa position dans la phrase. C'est ici qu'une bonne compréhension de la structure des

`tei:name` devient essentielle.

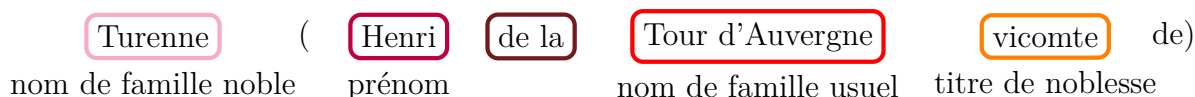


FIGURE 4.2 – Les différentes parties du `tei:name`

Il est donc nécessaire d'identifier les différentes structures possibles pour un `tei:name` contenant un nom de personne. Celles-ci sont visibles en 4.3. Pour les besoins de l'algorithme, deux structures principales ont été retenues : les `tei:name` contenant des parenthèses et eux n'en contenant pas<sup>6</sup>. Comme nous le verrons, d'autres subdivisions existent ensuite.

Les `tei:name` contenant des parenthèses ont une structure bien plus claire que ceux n'en contenant pas ; il est alors tout à fait possible d'identifier le rôle joué par les différents éléments de cette balise. L'organisation entre les différents types de noms est la même dans les exemples 4.4 et 4.5, qui contiennent tous les deux des parenthèses. Le nom le plus important, sous lequel une personne est connue au XIX<sup>ème</sup> s., est contenu entre parenthèses et au début de l'entrée, parfois en majuscules. Le ou les noms suivants sont, eux, à l'intérieur des parenthèses. Le premier nom au sein des parenthèses est toujours le prénom ; ensuite viennent des informations complémentaires. Ces deux exemples, cependant, forment deux sous-catégories dans le groupe des noms sans parenthèses :

- Le premier exemple (4.4) contient uniquement des noms de personnes nobles. Dans ce cas, le nom hors parenthèses correspond au nom de famille noble ; le nom de famille usuel d'une personne est contenu à l'intérieur des parenthèses. Ce détail est important, parce que les personnes sont souvent référencées sur *Wikidata* par le nom de famille usuel. Lancer une recherche sur l'API avec les noms de famille usuel et noble à la fois ne retournera pas toujours de réponse, alors que ne rechercher qu'un des deux noms peut permettre d'obtenir le bon résultat. C'est pourquoi, dans ces cas, les deux noms sont distingués : le nom de famille usuel est associé avec la clé `lname` du dictionnaire, et le nom noble lié au `nobname_sts`.
- Les personnes dans le second exemple (4.5) ne sont pas nobles, ce qui simplifie beaucoup l'extraction d'information. Il suffit de stocker la partie entre parenthèses hors du `lname` ; le prénom entre parenthèses, identifié grâce à l'algorithme d'extraction des prénoms, est stocké dans le `fname`. Parfois, d'autres informations sont contenues dans les parenthèses ; si elles sont significatives (comme dans le cas d'Alexandre Dumas père, où il est important de faire la distinction d'avec son fils), elles sont extraites. Sinon, ces informations supplémentaires sont abandonnées.

Les `tei:name` ne contenant pas de parenthèses (exemple 4.6) ne peuvent être traités avec la même granularité, puisqu'ils ne présentent pas de caractères récurrents pour faire

6. Comme cela étant déjà visible dans la figure 4.1

```

1 <name type="author">Turenne (Henri de la Tour d'Auvergne vicomte
  ↪ de)</name>
2 <name type="author">HUMBOLDT (Alexandre baron de)</name>
3 <name type="author">Tascher (Pierre-Jean-Alexandre-Jacquemin comte
  ↪ Imbert de)</name>

```

Code source 4.4 – Trois exemples de `tei:name` avec titres de noblesse

```

1 <name type="author">Viardot (Pauline)</name>
2 <name type="author">Verdi (Giuseppe)</name>
3 <name type="author">Sobieski (Thérèse-Cunégonde)</name>
4 <name type="author">CAUCHY (le bon Alex.)</name>
5 <name type="author">DUMAS (Alex. père)</name>

```

Code source 4.5 – Cinq exemples de `tei:name` sans titres de noblesse

la différence entre leurs différentes parties. Il est donc impossible de s'appuyer de façon récurrente sur la position des différents noms dans la phrase pour déduire leur signification. Cependant, il faut également remarquer que les éléments ne contenant pas de parenthèses contiennent en général bien moins d'informations que ceux qui en contiennent. Chercher à tout prix à identifier une structure à ces entrées afin de pouvoir avoir un dictionnaire complet à partir duquel lancer plusieurs requêtes n'a pas non plus nécessairement d'intérêt.

```

1 <name type="author">Sophie</name>
2 <name type="author">Henri VI</name>
3 <name type="author">DUPUIS</name>

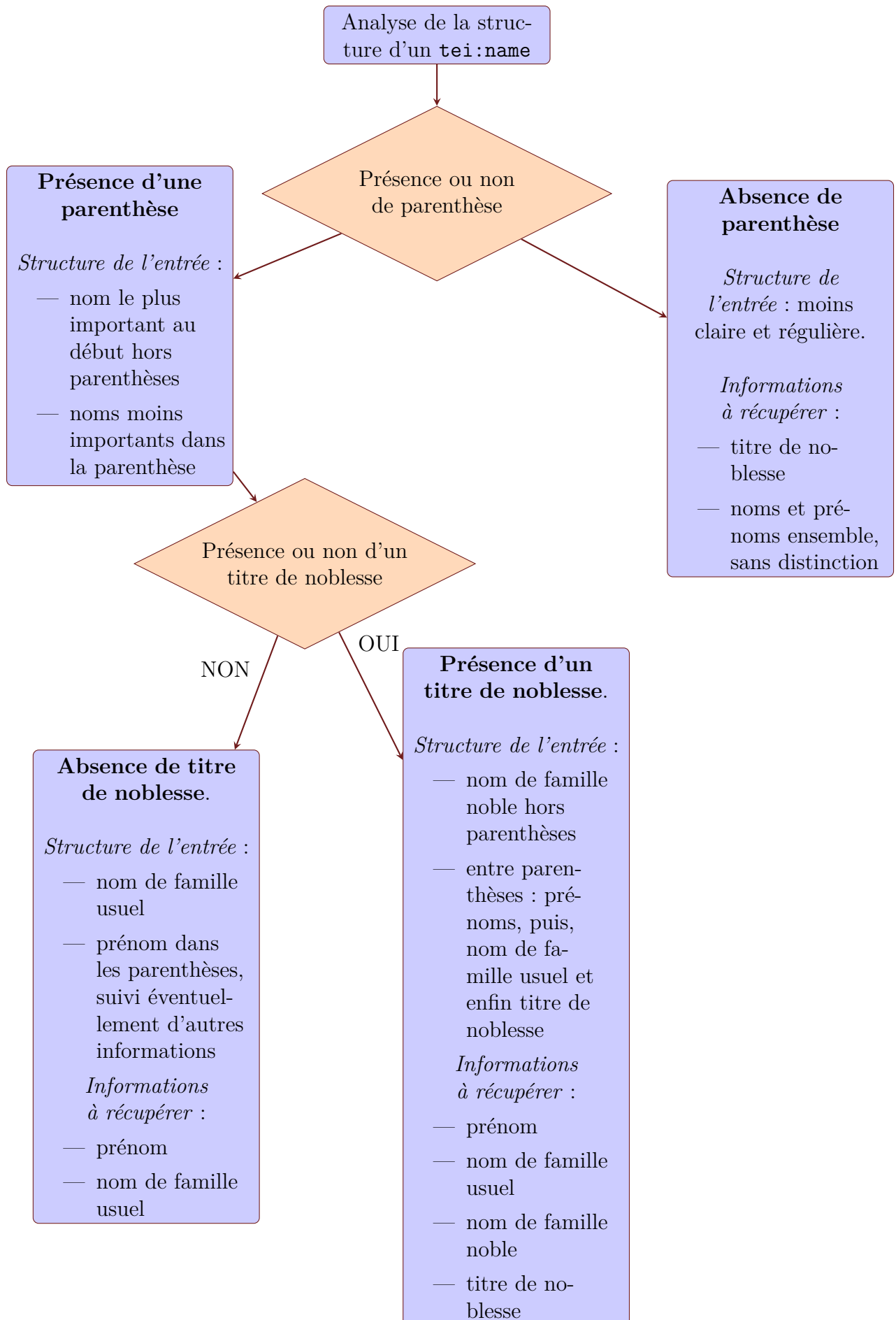
```

Code source 4.6 – Trois exemples de `tei:name` sans parenthèses

### 4.3.2 Identifier et extraire les informations nominatives

C'est sur les différents critères présentés dans la partie précédente et résumés sous forme graphique dans la figure 4.3, que se base l'extraction d'informations nominatives du `tei:name`.

Dans un premier temps, l'algorithme cherche à extraire des informations d'un `tei:name` contenant des parenthèses (étape représentée dans le graphique : 4.4). Il y a alors quatre données à identifier et à extraire : le nom de famille noble, le nom de famille usuel, le prénom et le titre de noblesse. L'algorithme commence par extraire des informations du texte contenu entre parenthèses avant de passer au texte hors parenthèses : il identifie un titre de noblesse avant d'extraire un prénom et éventuellement un nom de famille ; pour finir, le nom hors parenthèses est extrait.

FIGURE 4.3 – Différentes structures possibles pour un `tei:name`

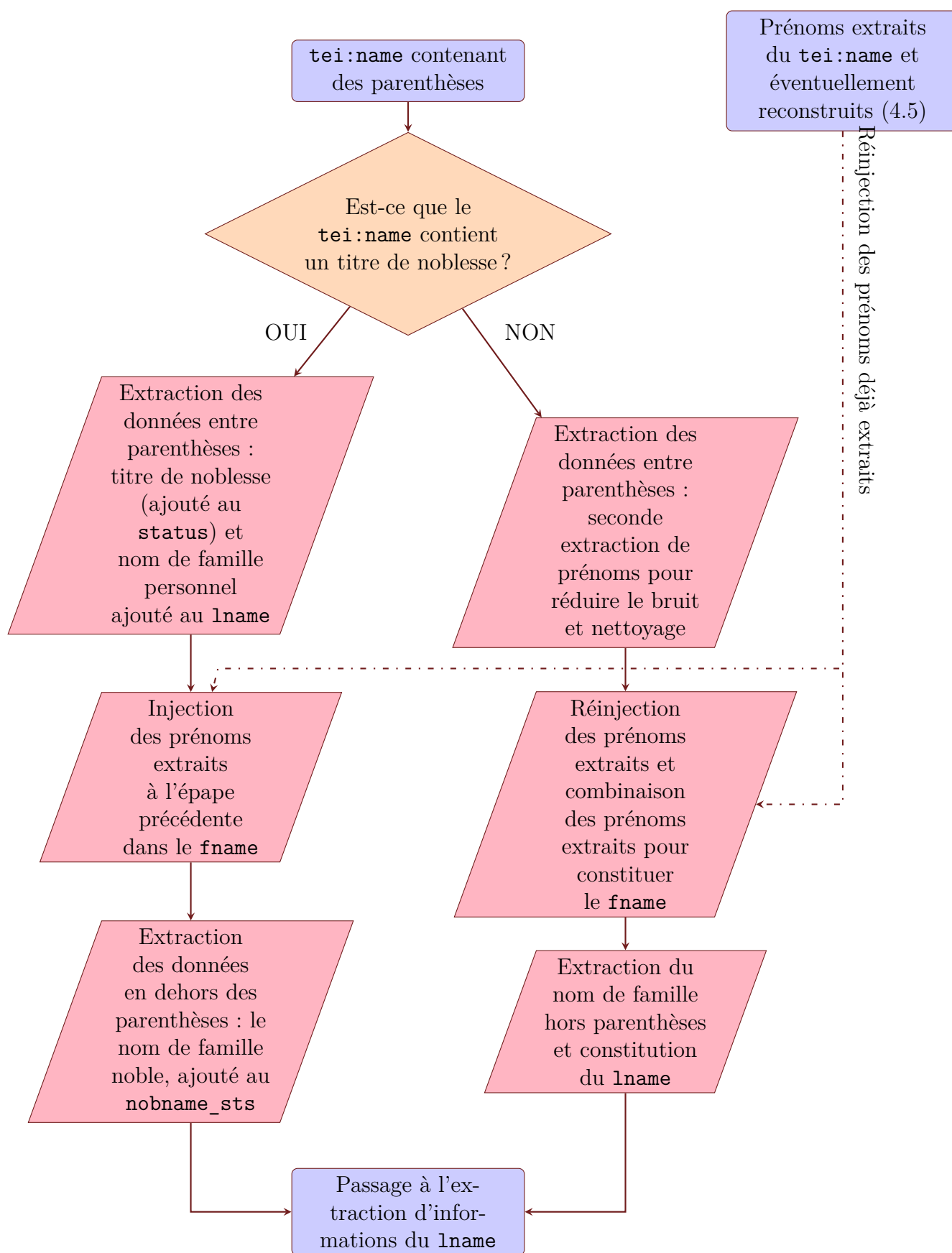
L'extraction du titre de noblesse est l'étape la plus simple, puisqu'elle repose sur une table de conversion (7) : l'algorithme cherche à identifier un titre de noblesse parmi les clés de ce dictionnaire. Si des titres sont trouvés, alors leur version normalisée figurant en valeur de ce dictionnaire est ajoutée aux données extraites. Il est à noter que, comme tous les titres de noblesse les plus importants permettent d'obtenir des résultats sur *Wikidata*, aucune valeur n'est associée aux titres les moins importants ; dans ce cas, aucun titre n'est extrait, puisqu'il risque d'empêcher l'obtention d'un résultat. Cependant, même des titres de noblesse n'ayant pas d'équivalent, et qui ne seront donc pas extraits, figurent sur cette table de conversion : la table permet d'identifier la présence d'un titre de noblesse, ce qui est nécessaire pour extraire les autres informations.

Si un titre de noblesse a été trouvé, un processus éliminatoire est engagé afin d'identifier le prénom et le nom de famille usuels contenus entre parenthèses. Comme cela apparaît dans les exemples plus haut, le texte entre parenthèses peut contenir du bruit : des mots comme « du, le... » qui n'aident pas à obtenir de résultats sur *Wikidata*, ou encore des noms communs et des attributs des personnes. Il faut détecter tous les noms propres, situer leur rôle et les extraire. Il peut être tentant d'extraire seulement les noms propres à l'aide d'expressions régulières – c'est à dire, d'identifier tous les mots commençant par une majuscule – et d'abandonner le reste. Cependant, certains noms communs peuvent être écrits avec une capitale ; en suivant cette méthode, il y aurait donc un risque d'inclure du bruit dans les données utilisées pour lancer les requêtes. L'algorithme lance donc une série de simplifications du contenu entre parenthèses : il commence d'abord par supprimer les données déjà extraites : le prénom ou son abréviation identifiées et le titre de noblesse ; ensuite, il supprime l'ensemble des termes « auxiliaires » (« le, la, dit... »), qu'ils aient des capitales ou non ; enfin, tous les noms communs restants sont supprimés. Il ne reste alors du texte entre parenthèses qu'un ou plusieurs noms commençant par une majuscule et considérés comme des noms propres ; quelle valeur donner à ce texte ? Étant donné que le prénom a déjà été identifié, de même que le titre de noblesse, le texte entre parenthèses ne peut être qu'un nom de famille ; on considère alors que le texte entre parenthèses est un de nom de famille usuel ; le texte hors parenthèses est alors le nom de famille noble.

Si un `tei:name` contient des parenthèses mais pas de titre de noblesse, le processus est analogue à celui décrit ci-dessus – mais plus simple, puisque le texte entre parenthèses a tendance à être bien plus complexe pour les personnes nobles. Le bruit est successivement détecté et supprimé à l'aide d'expressions régulières ; le prénom déjà extrait est supprimé. À ce stade, une deuxième extraction de prénoms a lieu : les prénoms étant parfois complexes et écrits d'une manière variable (avec ou sans parenthèses...), il arrive que ceux-ci ne soient pas entièrement extraits. Cette seconde extraction pose une légère difficulté, puisque l'on dispose alors de deux noms propres ayant la valeur d'un prénom. Il faut déterminer la relation entre ces deux « prénoms » : les mettre simplement bout-à-bout risquerait de créer des noms et prénoms dans le désordre ; dans ces cas, le moteur de

recherche de *Wikidata* ne retournerait pas de résultat. Par conséquent, c'est en fonction de la position relative des deux « prénoms » que le prénom final est reconstruit : celui qui se situe au début du texte entre parenthèse est positionné en premier dans le prénom reconstruit. Pour finir, c'est le texte hors parenthèses qui est extrait, soit le nom de famille. En l'absence de titre de noblesse, il n'est pas nécessaire de faire la différence entre différents noms de famille ; le nom hors parenthèses ne peut être qu'un nom de famille usuel.

Si un nom ne contient pas de parenthèses, le processus d'extraction est à la fois plus simple et moins qualitatif : il est tout simplement impossible de classer les différentes données extraites par détection de motifs. Dans ce cas là, le bruit supprimé du `tei:name` ; l'objectif est de pouvoir lancer des recherches sur *Wikidata* avec les données les plus propres possibles, afin de maximiser les possibilités d'obtenir un résultat correct.

FIGURE 4.4 – Extraction d'informations d'un `tei:name` contenant des parenthèses



### 4.3.3 Reconstruire un prénom complet à partir de son abréviation

Cette phase essentielle vise à corriger un problème dans les données : le nom d'une personne est souvent abrégé. Cela pose un problème pour l'utilisation ultérieure de l'API : le prénom est une information essentielle pour identifier l'alignement avec la bonne entité sur *Wikidata*. Sans ce prénom, le risque est grand que le premier résultat obtenu à l'aide du moteur de recherche de *Wikidata* ne soit pas celui de la personne décrite dans le `tei:name`, mais celui d'une autre membre de la même famille. De plus, si une requête est lancée avec un prénom abrégé (les initiales d'une personne, par exemple), elle ne retourne pas de résultat<sup>7</sup>. Le choix a donc été fait, pour chaque entrée, d'essayer de reconstruire un prénom complet à partir de son abréviation. Comme cela a été montré plus haut, l'extraction du prénom est également essentielle pour identifier les autres informations du `tei:name` par élimination. La reconstruction des noms a lieu en trois étapes :

- D'abord, un prénom, abrégé ou complet, est extrait.
- Ensuite, le type de prénom est identifié. Pour cet algorithme, trois types de prénoms existent : les prénoms complets (composés ou non), les prénoms abrégés non composés et les prénoms abrégés composés. Le choix de faire une distinction entre les deux catégories s'explique par la différence dans le processus de reconstruction
- Enfin, l'algorithme cherche à reconstruire le prénom à l'aide de tables de conversions.

La première question qui se pose est : qu'est-ce qui constitue un prénom ? Celui-ci correspond à un motif (un ou plusieurs noms propres), mais il est identifié par sa position dans la phase : le prénom est la première information contenue entre parenthèses – cela vaut autant pour les noms de personnes nobles que pour les autres. Pour être certain de cibler uniquement le prénom, cette étape n'a donc lieu que dans les `tei:name` contenant des parenthèses. Bien qu'un motif « prénom » puisse être identifié et que la position du prénom soit fixe, il est en fait assez difficile de modéliser ce motif de façon précise. Il faut définir un ou plusieurs motifs qui correspondent à tous les cas de figures, à toutes les formes que peuvent prendre un prénom. Voici quelques exemples pour mieux cibler cette difficulté (les prénoms sont en gras) :

- « Sobieski (**Thérèse-Cunégonde**) » : un cas classique : prénom composé où les deux noms sont complets et séparés par des tirets.
- « Zimmermann (**P.-J.-G.**) » : une autre forme simple : un nom composé où seules restent les initiales, écrites en majuscules, terminées par des points et séparées par des tirets.
- « DUBOIS-FONTANELLE (**J. Gaspard**) » : une forme un peu plus complexe : le nom est partiellement abrégé, avec une initiale et un nom complet ; cependant, il

---

7. Par exemple, voici une recherche sur Jean-Jacques Ampère. Si l'on recherche « J J Ampère », les résultats obtenus sont incorrects. En recherchant « Jean Jacques Ampère », le résultat obtenu est le bon.

n’y a pas de tiret qui sépare les deux prénoms, et il faut donc que les deux parties soient détectées toutes les deux dans le même prénom (sans quoi, il existe un risque de se retrouver avec des prénoms partiels).

- « DESCHAMPS (**Jh Fr. L.**) » : pour finir, une forme qui mélange les complexités : trois prénoms abrégés, mais qui ne contiennent pas qu’une initiale ; il ne sont pas tous terminés par des points et ne sont pas séparés par des tirets.

Cette variété tient autant du bruit dans les données que de différentes formes de notation dans le corpus. Elle complexifie en tout cas la définition, d’un point de vue formel, de ce qui constitue un prénom ; il est notamment difficile d’extraire l’intégralité d’un prénom composé, puisqu’il faut faire en sorte que les différents noms ou initiales soient reconnus comme faisant partie du même nom, qu’ils soient séparés par des tirets ou non. Dans l’abstrait.

- Un nom composé abrégé est une série de majuscules suivies d’autres lettres et/ou de points ; il doit contenir au moins un point et un tiret (un motif plus précis risque de ne pas identifier tous les cas de figure ; un motif plus généraliste risque d’inclure trop de bruit).
- Un nom simple abrégé correspond à une majuscule, suivie ou non d’autres lettres et terminée par un point.
- Un nom complet est une série de majuscules suivies de minuscules, séparées ou non par des tirets.

Identifier ces motifs « dans l’abstrait » est une étape importante, qui permet de mettre en place une implémentation technique. Cependant, il est impossible de construire un motif unique correspondant à tous les cas de figures. La décision a donc été d’accepter cette variété et ce bruit et de créer une série de motifs sous la forme d’expressions régulières, qui permettent de cibler différents cas de figure.<sup>8</sup> À partir d’une connaissance du corpus et de sa structure, il est donc possible d’identifier tous les prénoms dans le texte, et d’automatiser cette identification à l’aide de langages de programmation détectant des motifs.

Une fois les prénoms identifiés, il s’agit de reconstruire ceux qui sont abrégés. Comme pour le reste du processus d’extraction d’informations, cette étape s’appuie fortement sur l’usage de tables de conversion. Le principe général est simple : des tables de conversion sous la forme de dictionnaires contiennent en clés des formes abrégées de noms, et en valeurs les formes complètes. Si un nom extrait se retrouve en clé de ces tables, alors il est remplacé par sa version complète. Malheureusement, l’implémentation technique n’est pas aussi facile. En regardant les tables de conversions utilisées (8, 9), il apparaît que celles-ci sont assez courtes, et ne contiennent pas une grande variété de noms. Le

---

8. Les trois fonctions permettant l’extraction de prénoms sont visibles ici ; 10 permet d’identifier les noms composés abrégés, 11 identifie des noms abrégés non composés et 12 se consacre aux noms complets

problème avec l’usage de tables de conversion, c’est d’abord qu’elles doivent être faites à la main, et qu’il n’est pas possible de prendre toutes les abréviations possibles en compte ; mais le véritable problème est qu’une seule abréviation peut faire référence à plusieurs noms. C’est malheureusement le cas des abréviations les plus banales : « J. » peut à la fois signifier « Jean », « Jeanne », « Joséphique »... Le choix a donc été fait de ne retenir que les abréviations « évidentes » pour éviter d’introduire du bruit supplémentaire. De plus, quand un nom est modifié par l’algorithme de reconstruction, cela est indiqué à l’aide de la clé `rebuilt` du dictionnaire utilisé pour lancer des requêtes (4.2). Du fait de cette relative pauvreté des données, l’approche est ici totalement prédictive : il est impossible d’être certain d’obtenir le bon nom complet à partir de son abréviation ; il est seulement possible de prédire que le prénom reconstruit sera conforme au vrai prénom (tel qu’il est écrit sur *Wikidata*). Par son fonctionnement modulaire, qui s’adapte aux différents cas de figure, l’algorithme de reconstitution des prénoms cherche à maximiser cette certitude.

L’algorithme de reconstruction des prénoms a été conçu pour chercher à palier à cette relative pauvreté des données de conversion, surtout dans le cas des noms composés abrégés. Ceux-ci posent en effet un problème au delà de leur simple extraction : presque toutes les combinaisons de prénoms sont possibles. L’algorithme cherche donc à reconstruire un prénom composé en plusieurs étapes, d’abord en utilisant des tables de conversion dédiées aux noms composés, puis en utilisant des tables ne comportant que des noms composés. Pour plus de clarté, ce processus sera décrit à partir de l’exemple de « J.-P.-Ch. », à partir duquel il faut arriver à « Jean-Pierre-Charles ». La reconstruction d’un prénom composé commence par une tentative d’alignement complet un prénom composé abrégé présent dans les tables de conversion : un remplacement ne sera fait que si prénom correspond parfaitement aux données présentes dans la table visible en (8). Si il y a alignement complet, alors l’extraction du prénom s’arrête ici. Dans notre exemple, ce n’est pas le cas. Il faut tout de même tenter de reconstruire un prénom, ne serait-ce que partiel. Pour minimiser le risque d’erreur, cette reconstruction partielle se fait en deux temps. Dans un premier temps, l’algorithme cherche à associer une partie des initiales extraites avec des clés de la table consacrée aux noms composés. Dans l’exemple, « J.-P. » serait transformé en « Jean-Pierre », puisque « j p » se trouve dans la table de conversion. Si l’intégralité du nom abrégé est reconstruite à l’issue de cette étape, l’algorithme prend fin. Dans notre exemple, il reste à transformer le « Ch. » en « Charles ». L’algorithme essaye donc ensuite d’aligner les initiales restantes avec des clés de la table consacrée aux noms non composés (9). « ch » étant présent dans cette table, le dernier nom abrégé est remplacé par un nom complet : un nom a donc été entièrement reconstruit à partir de son abréviation. Tout au long de ce processus, un système permet de suivre quelles initiales ont été remplacées et lesquelles restent à être traitées, pour éviter de traduire deux fois une seule initiale. À la fin de ce processus, toutes les initiales qui n’ont pas été remplacées sont supprimées : si les abréviations étaient conservées, cela compromettrait

l'obtention de résultats sur l'API. À l'inverse, si seulement une partie des prénoms d'une personne sont extraits, mais que les abréviations sont supprimées, la possibilité d'obtenir des résultats sur l'API augmente.

L'alignement d'un prénom abrégé non-composé avec son équivalent complet est bien moins complexe : l'abréviation d'un prénom non-composé étant plus simple, elle ne permet pas de concevoir d'algorithmes réactifs qui puissent adopter différents comportements en fonction des résultats obtenus. Une fois un prénom abrégé extrait, l'algorithme cherche à l'aligner avec son équivalent complet présent dans la table 9. Si un prénom complet est reconstruit, celui-ci est utilisé pour lancer des requêtes sur l'API de *Wikidata*. Sinon, l'abréviation est supprimée pour ne pas mettre en danger l'obtention de résultats.

Parmi les trois types de prénoms détectés – prénom composé abrégé, prénom non-composé abrégé et prénom complet –, il ne reste donc que les prénoms complets. En théorie, ceux-ci ne demandent pas de retraitement. Cependant, des erreurs peuvent subsister dans les données utilisées. Il est possible que les « marqueurs » d'initiales et d'abréviations (le point à la fin d'un mot commençant par une majuscule) soient absents des documents sources ; le point étant un caractère peu visible, il est également possible que celui-ci soit présent dans les documents d'origine mais qu'il ne soit pas reconnu lors de l'extraction du texte par océrisation. Quelle que soit l'origine du problème, cela mène à des situations comme cela : `<name type="author">MALET (Cl-François)</name>`. Ici, « Cl » est une abréviation ; la structure du nom est cependant celle d'un nom propre composé. La détection de motifs ne pouvant s'appuyer que sur la forme du texte, et non sur sa signification, ce nom est reconnu comme un nom complet. Lorsqu'un prénom est identifié comme étant non abrégé, une étape supplémentaire de vérification a donc lieu : l'algorithme cherche à associer le ou les prénoms avec des prénoms abrégés présents dans les tables de conversions. Si ils se trouvent effectivement dans des tables, alors ils sont remplacés par leur équivalent complet. Sinon, par élimination, les prénoms sont considérés comme n'étant pas abrégés et sont utilisés pour lancer des requêtes sur l'API.

Ce processus de reconstitution des prénoms permet d'augmenter la qualité des données en cherchant à pallier à leur plus grand écueil. Bien que ce processus soit totalement prédictif et qu'il vienne avec un certain risque d'erreur, son rôle est essentiel dans l'obtention de résultats valides sur le moteur de recherche de l'API *Wikidata*, surtout en l'absence d'informations biographiques. Cela apparaît dans les tests menés, et dont le résultat est visible en annexes (1). Lorsqu'une requête est lancée seulement sur le prénom et le nom de famille usuel, l'identifiant *Wikidata* récupéré est valide dans 48% des cas en reconstituant les prénoms. Si des requêtes sont lancées sur les noms de famille et les prénoms non reconstitués uniquement, ce taux de réussite tombe à 42%. Cette différence est non négligeable lorsque l'on travaille sur 82000 entrées différentes à requêter. Ainsi, en utilisant une approche modulable et adaptée au corpus traité, il est possible de pallier à certaines absences dans les données disponibles. S'il n'y a pas toujours assez d'informa-

tions pour s'assurer que l'identifiant récupéré sur *Wikidata* sera le bon, reconstituer les prénoms permet d'augmenter leur précision des données, et donc d'obtenir de meilleurs résultats. Comme cela a été dit plus tôt, cette approche est totalement prédictive, et il est difficile de garantir qu'un prénom reconstruit sera correct. Cette incertitude ne doit pas être oubliée dans la constitution du programme chargé de lancer les requêtes en tant que tel : toute opération d'extraction et de reconstitution des données risque d'introduire du bruit. Il est donc important de lancer les recherches sur l'API de manière, là encore, modulaire, en multipliant le nombre de recherches pour chercher à pallier l'imperfection des données.

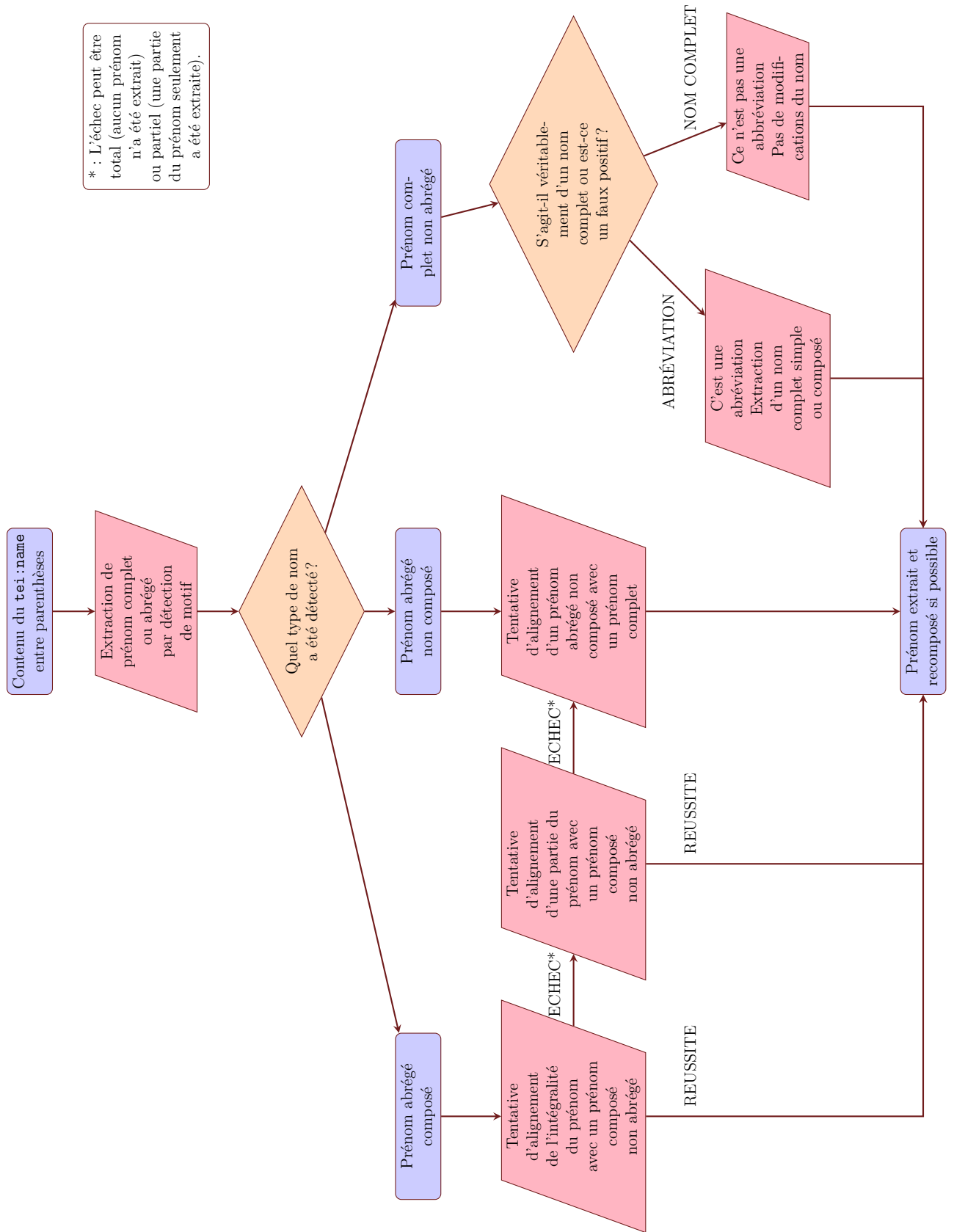


FIGURE 4.5 – Représentation graphique du processus d'extraction et de reconstitution d'un nom à partir de son abréviation

## 4.4 Extraire des informations biographiques du `tei:trait`

Si les informations nominatives sont déterminantes pour récupérer un identifiant valide, elles sont souvent incomplètes ; la manière dont les informations y figurent ne correspond pas non plus toujours à ce qui est disponible sur *Wikidata*. En plus du processus d'extraction d'informations nominatives issu du `tei:name` présenté ci-dessus, des informations biographiques sont extraites dans le `tei:trait`. Pour cette étape encore, l'extraction s'appuie sur une détection de motifs augmentée par l'usage de tables de conversion. Deux informations sont récupérées : les dates de vie et de mort – information cruciale car présentant un risque d'erreur très faible – et la fonction occupée par une personne – ce qui permet d'augmenter la certitude que l'identifiant *Wikidata* récupéré pour une entrée de catalogue est valide.

### 4.4.1 Identifier les dates de vie et de mort d'une personne

Comparativement à l'extraction d'informations nominatives, l'extraction du nom d'une personne par détection de motifs semble assez simple. Là où un nom peut prendre de nombreuses formes et structures différentes, l'expression régulière correspondant à une date correspond à une suite de quatre chiffres, soit :

`\d{4}`

S'il est facile de détecter une date, il est moins évident d'être certain que l'année extraite correspond à la naissance ou à la mort d'une personne. En effet, d'autres dates sont souvent présentes dans le `tei:trait`, comme cela apparaît dans l'exemple ci-dessous (4.7). Se contenter d'une extraction de toutes les dates du `tei:trait` risquerait de faire exploser la quantité de bruit dans les données utilisées pour lancer des recherches sur l'API *Wikidata*. Comme avec l'extraction d'informations du `tei:name`, détecter des motifs dans le texte ne permet pas d'attribuer une signification à ces motifs – il ne suffit pas de repérer une date pour que cette date corresponde à la naissance ou à la mort d'une personne.

Pour comprendre le sens des dates repérées dans le texte, il est encore une fois nécessaire de replacer cette date dans son contexte en s'appuyant sur la nature semi-structurée des documents. En effet, les `tei:trait` ayant tous une structure équivalente et utilisant un vocabulaire limité, il est possible de distinguer les dates pertinentes des autres. Les dates de naissance et de mort sont toujours des compléments dans les phrases, comme cela apparaît dans les différents exemples présentés ci-dessus (4.7). Dans l'exemple 5, par exemple, la date de naissance est un complément de « né » et la date de décès de « fusillé ». C'est donc en identifiant ce à quoi une date se rapporte qu'il est possible de distinguer les dates pertinentes des autres. En général, une date naissance est précédée par « né à, née à ». Le vocabulaire utilisé avec la date de décès d'une personne, quoique plus coloré, est la encore assez limité : « décapité.e, assassiné.e, fusillé.e, guillotiné.e, tué.e ».

```

1 <!-- exemple 1 -->
2 <trait>
3   <p>maire (en 1848) du XIIe arrondissement, médecin, sous-directeur
   ↪ de l'hôpital de la Salpêtrière.</p>
4 </trait>
5
6 <!-- exemple 2 -->
7 <trait>
8   <p>célèbre maréchal de France, vainqueur d'Alger en 1830, né en
   ↪ 1773, mort en 1846</p>
9 </trait>
10
11 <!-- exemple 3 -->
12 <trait>
13   <p>brave général de l'armée de l'Est pendant la guerre de 1870, né à
   ↪ Sarreguemines en 1840, mort en 1876</p>
14 </trait>
15
16 <!-- exemple 4 -->
17 <trait>
18   <p>célèbre révolutionnaire, membre de la Commune de Paris, né en
   ↪ 1838, tué à Rueil le 3 avril 1871.</p>
19 </trait>
20
21 <!-- exemple 5 -->
22 <trait>
23   <p>célèbre jurisconsulte, ministre et sénateur, un des ôtages de la
   ↪ Commune, né à Valence (Drôme) en 1804, fusillé avec l'archevêque
   ↪ de Paris le 27 mai 1871</p>
24 </trait>

```

Code source 4.7 – Exemples de `tei:trait` contenant des dates autres que les dates de naissance et de mort

Ce vocabulaire étant relativement limité, il n'est pas nécessaire d'utiliser de tables de conversions. Dans le motif permettant d'extraire une date, il suffit simplement de prendre en compte le contexte de celle-ci. Le motif à rechercher correspond donc à :

Mot indiquant une naissance ou un décès -- autres informations (lieu de naissance...) -- date

L'implémentation de ce motif sous la forme d'expressions régulières correspond donc à :

pour les dates de naissance : `(^\s|,|)[Nn](.\s|ée?)+?(?=\d{4})\d{4}`  
pour les dates de décès : `(^\s|,|)((M|m)|[Mm](\s|orte?))+?(?=\d{4})\d{4}`



et : `(^\s|,|)([Dd]écap|[Aa]ssa|[Tt]uée?|[Ff]usi|[Gg]uil).+?(?=\d{4})\d{4}`

Si ce motif est détecté, il est certain que la date est une date de naissance ou de décès. Il suffit alors d'extraire la date de son contexte. Ainsi, sur des documents semi-structurés, il est possible, à l'aide de détection de motifs, de resituer des éléments dans leur contexte afin de produire des données propres et réutilisables.

#### 4.4.2 Identifier l'occupation d'une personne

Identifier et extraire l'occupation d'une personne s'appuie également sur de la détection de motifs, cette fois-ci augmentée de tables de conversion : il existe une grande variété dans les différents types de métiers occupés par une personne ; de plus, un métier tel qu'il est décrit dans un catalogue n'a pas nécessairement d'équivalent sur *Wikidata*. Il est donc nécessaire de normaliser les données extraites.

Le processus d'extraction d'informations à l'aide de tables est assez similaire à l'extraction de titres de noblesse, par exemple : pour chaque `tei:trait` présent dans les catalogues, l'algorithme vérifie si un terme correspond à une occupation à l'aide de la table 1. Si c'est le cas, l'occupation telle qu'elle est mentionnée dans les catalogues est remplacée par son équivalent normalisé. Celui-ci n'est pas seulement une traduction du terme tel qu'il est présent dans le `tei:trait` : les occupations ou fonctions sont remplacées par une version normalisée, qui permette d'obtenir des résultats sur le moteur de recherche de *Wikidata*. Par exemple, le terme « maréchal » sera remplacé par « marshall » (sa traduction anglaise), tandis que des grades militaires moins élevés (« capitaine », « commandant ») seront remplacés par le terme générique « militaire ». L'extraction d'informations n'étant pas viser à obtenir des données définitives, il n'est pas nécessaire de préserver les termes et les catégories utilisées dans les catalogues.

Cette extraction d'informations « à l'aveugle » – qui se soucie uniquement des mots utilisés et non du contexte dans lequel il le sont – présente deux problèmes, visibles dans les exemples ci-dessous (4.8).

- Dans le premier exemple, une personne est à la fois désignée comme étant « sculpteur » et « auteur ». À l'aide de la table de conversion, les deux termes sont extraits et normalisés, ce qui donne « sculptor writer »<sup>9</sup>. Non seulement la transformation est fautive (il ne s'agit pas d'un écrivain), mais elle rend inutile l'extraction d'informations : les deux termes se contredisent et ne permettront pas d'obtenir de résultats sur *Wikidata*. Ce problème est récurrent, puisque des personnes ont souvent écrit ou réalisé une œuvre sans qu'elles soient connues comme écrivaines – comme cela apparaît dans le deuxième exemple, où le personnage est général, mais a écrit des mémoires. Il faut donc déterminer laquelle de ces deux fonctions est à retenir pour faire des recherches sur l'API *Wikidata*.

---

9. « sculpteur écrivain »

- Au premier problème de l'extraction de plusieurs fonctions s'ajoute une deuxième difficulté : les occupations extraites ne se rattachent pas nécessairement à la personne nommée dans le `tei:name`. C'est le cas dans le troisième exemple : « empereur » est présent dans la table de conversion – contrairement à son équivalent féminin, qui ne permet pas souvent d'obtenir de résultats. La fonction extraite sera donc « empereur / emperor », ce qui fausse les données extraites : c'est la femme de l'empereur qui est décrite dans cette entrée. Il est donc nécessaire de reconnaître lorsque, dans le `tei:trait`, une fonction ne se rattache pas à l'auteur.ice du manuscrit.

```

1 <!-- exemple 1 -->
2 <trait>
3   <p>célèbre sculpteur, auteur de la statue de la Résistance, érigée à
   ↪ Dijon en 1875 et enlevée par ordre du gouvernement, né à Nuits
   ↪ (Côte-d'Or) en 1815, mort en 1876</p>
4 </trait>
5
6 <!-- exemple 2 -->
7 <trait>
8   <p>célèbre général auteur de Mémoires, qui ont eu un grand succès,
   ↪ n. 1782, m. 1854.</p>
9 </trait>
10
11 <!-- exemple 3 -->
12 <trait>
13   <p>impératrice des Romains, reine de Hongrie, fille de
   ↪ Charles-Quint, femme de l'empereur Maximilien II, morte en
   ↪ 1603</p>
14 </trait>

```

Code source 4.8 – Difficultés liées à l'extraction de fonctions du `tei:trait`

Les deux problèmes trouvent des solutions différentes. Le premier problème naît de la différence entre les informations présentes dans les catalogues et celles qui sont nécessaires pour l'extraction d'informations de *Wikidata*. Dans un texte fait pour être lu par des humains, il n'y a aucun problème à mettre bout à bout plusieurs occupations, fonctions ou faits notables à propos d'une personne. Un.e lecteur.ice est capable de comprendre la polysémie des mots (« auteur », par exemple) et la relation des différentes fonctions entre elles. Dans le cas d'un catalogue, qui est fait pour vendre, il est normal que les descriptions de personnes soient mélioratives : la multiplication des hauts faits à propos d'une personne met en avant son statut, et justifie donc la valeur du manuscrit vendu. Les techniques de traitement automatique utilisées ne comprennent pas la polysémie. Comme cela a été dit plus haut, le moteur de recherche de *Wikidata* fonctionne mieux avec des termes spécifiques ; il est moins efficace lorsque le nom d'une personne est recherché avec

plusieurs fonctions en même temps. Il est donc nécessaire de déterminer l'occupation principale d'une personne. Après l'extraction des fonctions, les différents termes extraits sont comparés entre eux afin de déterminer le terme à retenir. La comparaison se fait entre les mots qui sont cooccurents, c'est-à-dire qui se retrouvent régulièrement ensemble, comme par exemple « marshall » et « military ». Dans ce cas, c'est la fonction la plus importante et la plus spécifique (« marshall ») qui est conservée, et le terme générique est abandonné. Les termes très génériques – comme le « writer » qui est traduit d'« auteur » – sont presque systématiquement abandonnés si d'autres termes ont été extraits : un.e auteur.ice n'est un.e écrivain.e que si ce terme est le seul utilisé pour le ou la décrire. Sinon, la personne est autrice de quelque chose sans que cela soit sa fonction principale.

Le second problème naît de l'algorithme d'extraction d'informations lui-même : la récupération de fonctions ou de métiers est non-contextuelle : si un terme est présent, il est récupéré. Un post-traitement est donc nécessaire afin de déterminer si le terme se rapporte à l'auteur.ice du manuscrit, ou à une autre personne. L'algorithme cherche à voir si le terme est précédé des prépositions « du, de la, des », comme c'est le cas dans le troisième exemple ci-dessus (4.8). Dans ce cas, il est considéré que l'information récupérée n'est pas relative à l'auteur.ice, mais à une tierce personne et elle est donc abandonnée.

La récupération de données biographiques du `tei:trait` marque la fin du processus d'extraction d'informations ; à l'issue de celui-ci, l'ensemble des éléments récupérés sont stockés dans un dictionnaire structuré (4.2), ce qui permet d'attribuer à chaque terme une fonction, une signification spécifique. À partir d'un texte en langage naturel, il est donc possible d'extraire un certain nombre de paramètres à partir desquels lancer des recherches en plein texte. C'est grâce à la nature semi-structurée des entrées – qui utilisent un vocabulaire peu varié dans des phrases structurellement proches les unes des autres – qu'il est possible de faire une extraction d'informations avec un grand degré de précision en s'appuyant uniquement sur de la détection de motifs. Grâce à une bonne connaissance du corpus et à l'aide d'une approche modulaire, qui s'adapte à différents cas de figure, il est possible d'approximer une compréhension « humaine » du texte. L'extraction d'informations n'étant qu'un étape préalable, cette compréhension du texte est uniquement centrée sur les éléments qui sont pertinents pour le moteur de recherche de *Wikidata*. Cette étape d'extraction est cependant prédictive : elle peut inclure du bruit dans des données, les informations extraites peuvent être mal structurées, puisque leur fonction est déterminée principalement en fonction de leur place dans le texte. De ce fait, l'étape suivante se doit de répondre à cette part de bruit possible, à l'aide d'un algorithme qui soit ici encore modulable. L'attribution à chaque terme d'une valeur ou d'une fonction spécifique est mise à profit pour pouvoir construire de façon automatique des recherches intelligentes, qui s'adaptent autant aux données disponibles que à l'outil auquel elles sont destinées : un moteur de recherche.



## Chapitre 5

# Résoudre les entités par l'alignement avec des identifiants *Wikidata*

En utilisant les données extraites à l'étape précédente, des recherches en plein texte sont faites sur *Wikidata* de façon automatique, grâce à l'API développée par cette institution. *Wikidata* retourne plusieurs résultats ; il est impossible de vérifier quel résultat est pertinent parmi tous ceux qui sont obtenus, et seul le premier résultat est donc retenu. Ce résultat contient un identifiant *Wikidata*, à partir duquel des requêtes pourront être menées avec SPARQL afin de constituer un jeu de données sur les auteurs figurant dans les catalogues de vente de manuscrits. L'objectif de l'algorithme est alors de chercher à maximiser la probabilité que, avec les données produites à l'étape précédente, le premier résultat retourné par l'API soit correct. Pour savoir comment maximiser cette probabilité, il est nécessaire de mieux comprendre la fiabilité des différents paramètres extraits à l'étape précédente. Après avoir étudié l'impact des différents paramètres dans l'obtention d'un résultat valide et présenté le fonctionnement d'une API, la manière dont l'algorithme a été constitué sera présentée.

### 5.1 Quantifier l'incertitude : quels sont les paramètres qui permettent d'obtenir le bon résultat ?

Tout au long de la constitution de cet algorithme, des tests ont été menés pour maximiser sa qualité et la probabilité d'obtenir des résultats valides. Le premier test est essentiel à la constitution de l'algorithme : il mesure l'impact de chaque paramètre de recherche (dates de naissance et de décès, titres...) dans l'obtention d'un identifiant valide (1).

Ce test a été mené sur un jeu de 200 entrées sélectionnées dans différents catalogues, avec une répartition des différents types de `tei:name` similaire à celle du jeu de données complet ; de plus, la proportion d'entrées avec des `tei:trait` est la même dans le jeu de

test et dans l'ensemble du catalogue. Pour chacune de ces 200 entrées, des données sont extraites du `tei:name` et du `tei:trait`; ensuite, cinq requêtes sont lancées pour chaque entrée : d'abord une requête avec le nom de famille usuel et le prénom seulement ; ensuite, une requête avec ces deux informations et un paramètre en plus – nom de famille noble, titre de noblesse, dates, occupation. Les cinq requêtes sont lancées, peu importe si les différentes données sont disponibles pour une entrée ou non : l'intérêt est de pouvoir quantifier l'impact de chaque paramètre pour l'ensemble des entrées, et pas seulement pour celles où le paramètre est disponible. Par exemple, si les dates de naissance et de décès permettent une obtention du bon résultat dans 90% des cas mais qu'elles sont disponibles pour 5% entrées, leur impact global sur l'obtention d'un résultat valide est en fait très limité. Comme cela a été dit plus haut, ces cinq requêtes sont lancées deux fois : d'abord avec le nom et le prénom complet reconstitué si besoin, ensuite avec seulement les prénoms non-reconstitués, afin de vérifier l'utilité de la reconstitution (les résultats sont très positifs, avec une augmentation de performance allant jusqu'à +6%).

Ces tests permettent donc de connaître, au niveau de l'ensemble du jeu de données, les paramètres les plus fiables pour obtenir des résultats valides. Il en ressort que le paramètre le plus efficace est l'occupation, ou le métier d'une personne : celui-ci permet d'obtenir des résultats valides pour 53,1% des entrées (contre 48% en recherchant les noms et prénoms, et 42% en cherchant les noms et uniquement les prénoms non-reconstruits). Ce résultat peut sembler inattendu puisque les fonctions ne sont pas toutes relevées et qu'elles subissent un retraitement, ce qui pourrait créer du bruit additionnel dans les recherches. Cependant, il montre l'importance de la connaissance à la fois de son corpus et de *Wikidata* : c'est en s'adaptant à son moteur de recherche qu'une extraction efficace des fonctions a été rendue possible. Le deuxième paramètre le plus efficace est les dates de naissance et de mort (52,6% de résultats valides) : c'est une information qui n'est pas retraduite, et il n'y a pas de risque de création de bruit additionnel. Suivent les noms de famille noble (45,9% de 40,5%, respectivement avec et sans reconstitution des prénoms). Il est intéressant de remarquer que ce paramètre fait diminuer les performances des recherches (entre 1,5 et 2,1 points de moins qu'en utilisant uniquement prénom et nom de famille usuel). Cette chute s'explique probablement par le fait que, sur *Wikidata*, une personne est rarement nommée à la fois par son nom usuel et par son nom de famille noble. Dans certains cas cependant, rechercher le nom de famille noble à la place du nom usuel peut permettre d'obtenir des résultats qui n'auraient pas été trouvés autrement.

En quantifiant l'impact de chaque paramètre de recherche dans l'obtention d'un résultat, ce test permet de savoir à quelles données se fier et quelles données rechercher en premier sur l'API de *Wikidata*. Il permet donc de quantifier la certitude, ou l'incertitude, liée à chaque paramètre, et donc de développer un algorithme de recherche de façon plus intelligente. Mais avant de présenter cet algorithme, il peut être intéressant de revenir sur ce qu'est une API, et sur la manière de les utiliser dans un contexte de recherche.

## 5.2 Automatiser l'extraction de données depuis des sources externes : les API

### 5.2.1 Qu'est-ce qu'une API ?

Une API Web<sup>1</sup> est un protocole qui permet à un programme d'interagir avec un autre dans le cadre d'une architecture client-serveur. Cette architecture est à la base du Web, qui est formé d'un ensemble d'ordinateurs reliés en réseau. Certains ordinateurs sont des hôtes et servent à diffuser du contenu ou des données : ce sont les serveurs. D'autres ordinateurs sont des clients, c'est-à-dire qu'ils consomment ce qui est fourni par les serveurs<sup>2</sup>. Le client lance des requêtes qui sont traitées par le serveur, qui en retour renvoie une réponse. Ce réseau d'ordinateurs communique grâce à des protocoles, au cœur desquels se trouve le *HyperText Transfer Protocol* (HTTP)<sup>3</sup>. Le rôle de celui-ci est de définir comment un client et un serveur communiquent : il établit la sémantique des requêtes envoyées d'un client à un serveur, et des réponses que ce dernier renvoie en retour. Le protocole HTTP définit donc de façon générale comment deux machines communiquent : lorsqu'un client pose une question, le serveur saura toujours ce que la question signifie ; de la même manière, le client pourra toujours interpréter la réponse du serveur. Cependant, ce protocole ne définit pas ce qui peut être demandé par un client à un serveur ; il ne définit pas non plus précisément le modèle de données utilisé par la réponse du serveur. Des technologies et des protocoles supplémentaires sont donc nécessaires, pour définir le contenu des requêtes et des réponses, en plus de leur sémantique. C'est ici que les API prennent leur importance : elles définissent ce qui peut être demandé, le vocabulaire spécifique à utiliser par les clients et les formats de réponse qui seront fournis en retour. En plus de définir un protocole, elles sont des programmes complets côté serveur : lorsqu'elle reçoit une requête d'un client, l'API traite celle-ci, construit une réponse en interagissant avec serveur et la renvoie à un client. L'API est donc un programme, avec des fonctions qui traitent des données, et un protocole, avec un vocabulaire permettre l'interaction client-serveur. Contrairement au HTTP qui est générique et défini à l'échelle du Web, les API sont spécifiques à un projet, à une application (chaque application ayant des besoins différents) ou à un groupe d'applications<sup>4</sup>. Une API permet donc l'interaction entre deux types de programmes : le programme dit fournisseur (qui est soit l'API en tant que telle, soit un programme côté serveur) et le consommateur (c'est-à-dire, un programme côté

---

1. Il existe également des API qui ne sont pas faites pour le Web, et qui permettent par exemple la communication entre deux programmes présents sur un même ordinateur.

2. Dans les faits, une seule machine peut à la fois être client et serveur.

3. Roy Fielding, *Architectural Styles and the Design of Network-based Software Architectures*, Version : PDF edition, 1-column for viewing online, Thèse de doctorat, Irvine, University of California, 2000, URL : <https://www.ics.uci.edu/~fielding/pubs/dissertation/top.htm> (visité le 05/08/2022), p. 116-127.

4. Il existe cependant des standards dans l'architecture des API, comme le *Representational state transfer* (REST) , défini par R. Fielding (*Ibid.*, p. 76-106)

client)<sup>5</sup>. D'un point de vue de la conception d'applications, les API ont de nombreux avantages. Elles garantissent la sécurité des données : selon les principes REST<sup>6</sup>, le client n'accède jamais directement à une base de données dans le serveur, mais seulement à une représentation des données produite à la demande ; puisque les données enregistrées sont distinctes des données communiquées, il est également possible de modifier la base de données derrière l'application sans pour autant devoir modifier l'intégralité de l'API ; de la même manière, le programme client n'a pas besoin d'être mis à jour à chaque mise à jour côté serveur : il doit seulement être modifié si l'API est modifiée.

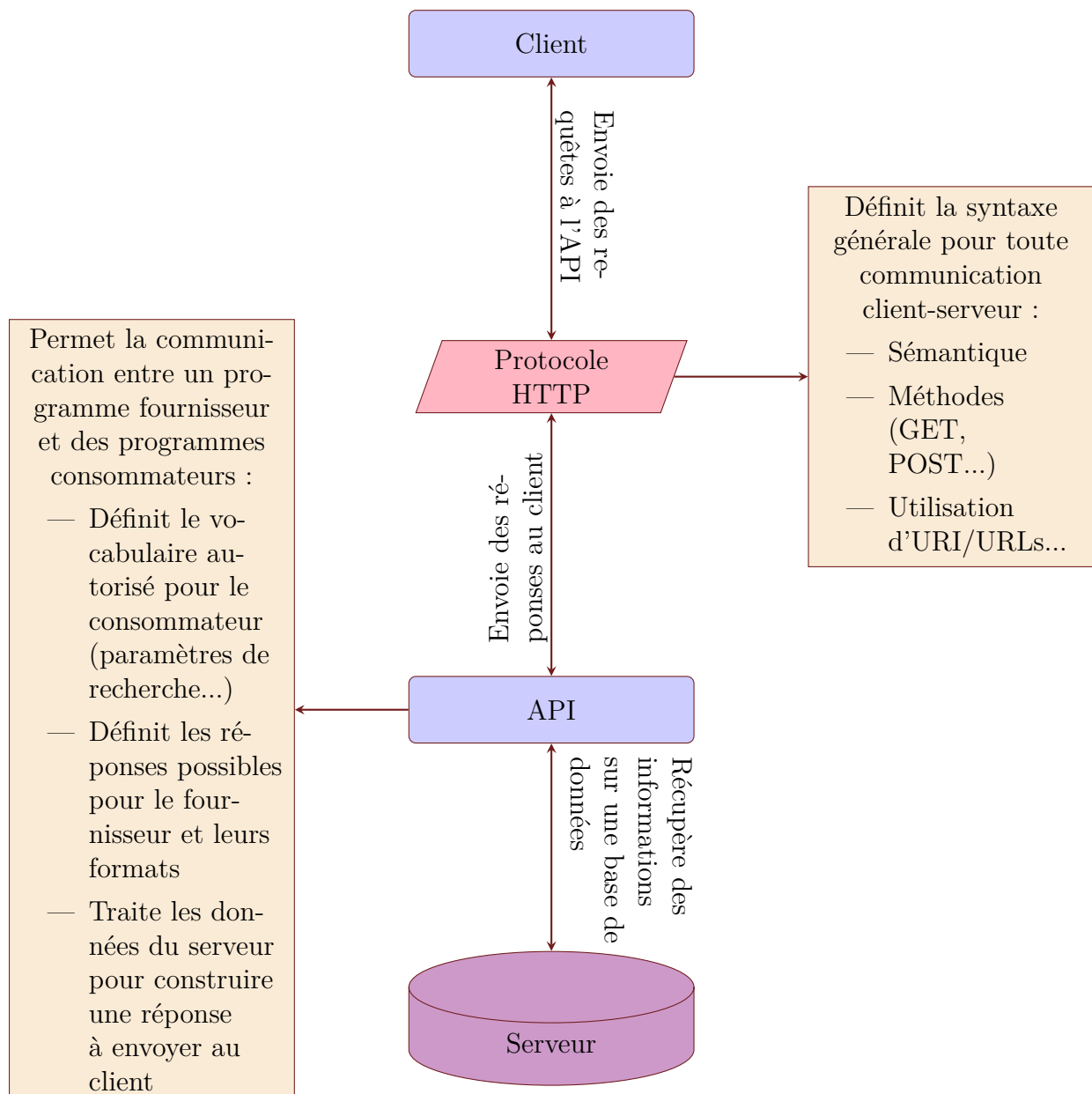


FIGURE 5.1 – L'interaction client-serveur au travers d'une API

5. Pour une représentation graphique de la communication avec une API dans une architecture client-serveur, voir la figure 5.1

6. *Ibid.*



Si l'élaboration d'une API est intéressant d'un point de vue de développement, d'architecture Web et d'interaction entre machines, l'utilisation d'API dans un contexte de recherche en humanités peut être très fructueuse. Une API permet la communication de données brutes (c'est-à-dire, d'informations « à l'état pur », sans mise en page ou interface graphique) dans des formats documentés. Consommer une API<sup>7</sup> permet donc d'obtenir de façon automatisée des données déjà structurées, dans un format qui soit compréhensible par une personne humaine autant que par une machine. L'utilisation d'une telle technologie est donc totalement adaptée à l'alignement avec *Wikidata*. Concrètement, en utilisant l'API de *Wikidata*, il est possible de lancer des recherches en plein texte sur le moteur de recherche de cette base de données de façon automatique. Même si les données obtenues doivent être ensuite corrigées, en écrivant un programme pour consommer une API, on s'épargne le travail de faire manuellement plus de 80000 recherches.

### 5.2.2 Consommer une API : le cas de l'application *MediaWiki*

*Wikidata* et les projets qui y sont affiliés (*Wikipedia*, *Wikimedia* et autres projets de la *Wikimedia Foundation*) peuvent être consultés en ligne, par des utilisateurs humains, au travers d'un site Web traditionnel. Il faut cependant distinguer le site Web des données : le site n'est qu'une interface de présentation. Les projets de la *Wikimedia Foundation* existent avant tout sous la forme d'un ensemble de bases de données (4), qui sont accessibles au public via le logiciel *open source* MediaWiki<sup>8</sup>. Ce logiciel permet la diffusion des données aux utilisateurs via une interface graphique traditionnelle ; mais il est également possible de récupérer et d'envoyer des données de la *Wikimedia Foundation* via l'API de MediaWiki.

Une API Web étant un équivalent d'une application Web conçue pour une machine, la manière d'accéder aux données est assez similaire à celle utilisée pour accéder à un site Web traditionnel : la requête passée par le client prend la forme d'un *Uniform Resource Locator* (URL)<sup>9</sup> analogue à celui d'une page Web traditionnelle. L'URL contient la plupart des informations nécessaires pour que l'API comprenne la requête (interprète ce que le client demande) et construise une réponse (dans notre cas, renvoie des données). Ci-dessous, deux URL correspondant à une recherche en plein texte pour « Virginia Woolf ».

Recherche sur le site *Wikidata* :

<https://www.wikidata.org/w/index.php?search=virginia+woolf>

---

7. C'est-à-dire récupérer des données

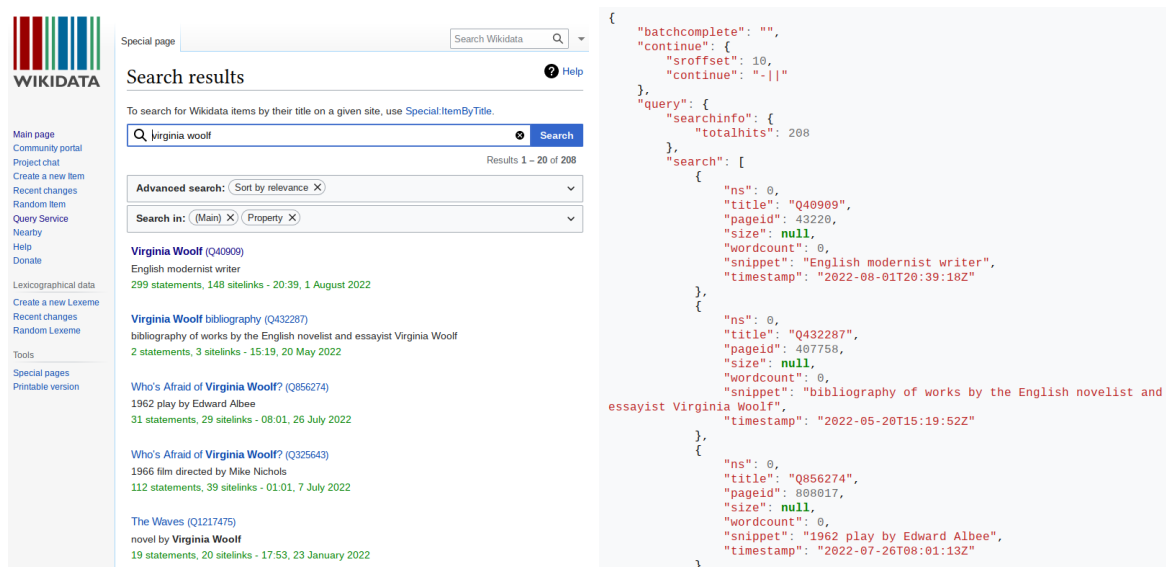
8. *Wikimedia Foundation\$Technology*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://en.wikipedia.org/wiki/Wikimedia\\_Foundation#Technology](https://en.wikipedia.org/wiki/Wikimedia_Foundation#Technology) (visité le 01/08/2022).

9. Dans les faits, dans le protocole HTTP les URL sont nécessaires à une requête, mais ne sont qu'une partie de la requête : elle est également composée d'un en-tête et d'une méthode, c'est-à-dire d'une action à faire avec cet URL (demander des données, en envoyer, en supprimer...).

Recherche sur l'API *Wikidata* : `https://www.wikidata.org/w/api.php?action=query&list=search&srsearch=virginia+woolf`

Les deux URL posent la même question ; cependant, le site « pour humain.e.s » est indiqué par `index.php?` ; `api.php?` indique que le deuxième URL renvoie à l'API et contient une suite de paramètres de recherches séparés par le caractère `&` :

- `action=query` : l'action effectuée par cette URL est une requête demandant des données au serveur.
- `list=search` : il s'agit d'une recherche plein texte.
- `srsearch=virginia+woolf` : rechercher les pages qui contiennent « Virginia Woolf ».



(a) Résultats retournés par le site traditionnel

(b) Résultats retournés par l'API

FIGURE 5.2 – Résultats retournés par *Wikidata* pour la recherche « Virginia Woolf »

Enfin, les données récupérées sont les mêmes, mais dans deux formats très différents (5.2). Dans l'image de gauche (5.2a), les résultats sont visibles avec une interface graphique complète, comme sur n'importe quel site Web. À droite cependant (5.2b), les résultats sont présentés sous la forme de données structurées au format JSON (format analogue à un dictionnaire). Ce résultat contient à la fois les résultats à proprement parler (en dessous de « Search ») et des informations descriptives. Par exemple, `query/searchinfo/totalhits` retourne le nombre de résultats pertinents, soit 208 pages.

Cette brève présentation de la manière de consommer concrètement une API indique également le protocole à suivre pour l'alignement des `tei:name` avec des entités *Wikidata* : il s'agit de concevoir un algorithme qui, à partir des données extraites du `tei:name` et du `tei:trait`, construise des URL pour lancer des recherches en plein texte sur *Wikidata*. Cela revient à associer au paramètre de l'URL `srsearch` la chaîne de caractère qui doit

être recherchée sur la base de données de *Wikidata*. Ensuite, l'algorithme doit transmettre à l'API ces , interpréter le JSON renvoyé par *Wikidata* et enregistrer ces réponses dans un fichier. Étant donné que l'alignement avec *Wikidata* prend plusieurs dizaines d'heures et demande des performances élevées, des fichiers de log annexes sont créés en cours d'utilisation pour optimiser l'algorithme.

## 5.3 Présentation générale

L'algorithme d'alignement avec *Wikidata* construit des URL et lance des recherches en plein texte sur l'API *Wikidata*. De par son fonctionnement, il multiplie les requêtes. Puisqu'il est impossible de vérifier en cours d'exécution si la réponse retournée par l'API contient un résultat valide, il est nécessaire d'organiser l'algorithme afin de maximiser la probabilité d'obtenir un résultat valide. C'est pourquoi les recherches sont faites en commençant par les celles qui ont la plus grande certitude d'obtenir des résultats valides et en terminant par les plus incertaines. Dès qu'une réponse est obtenue, l'algorithme s'arrête et l'entrée de catalogue suivante est traitée. Comme cela apparaît dans la figure 5.3, il a par conséquent un fonctionnement à plusieurs étapes.

La première étape est la même quelles que soient les informations contenues dans le `tei:name` : l'ensemble des paramètres<sup>10</sup> extraits lors de l'étape précédente sont mis bout à bout afin de lancer une recherche. Cela permet que la recherche soit la plus précise possible, afin de maximiser les possibilités d'obtenir un résultat valide.

Si la réponse de l'API ne contient pas de résultat, alors les données disponibles sont étudiées afin d'adopter le bon comportement. Comme cela a été dit, différents processus ont été utilisés pour extraire des données du `tei:name` et du `tei:trait` ; il est donc logique de ne pas traiter tous les dictionnaires de données de la même manière.

Si l'entrée concerne une personne noble (c'est-à-dire, si le dictionnaire de données contient un nom de famille noble), alors il y a un possible conflit dans les données extraites : si le nom de famille usuel est recherché en même temps que le nom de famille noble sur *Wikidata*, il arrive qu'aucun résultat ne soit retourné. Dans ce cas, une première requête est lancée en supprimant le nom de famille noble ; si aucun résultat n'est obtenu, celui-ci est rajouté et c'est le nom de famille usuel qui est enlevé. Pour finir, une autre recherche est lancée sans le prénom. Si à ce stade aucun résultat n'a été obtenu, alors un algorithme de recherches soustractives est lancé ; celui-ci est décrit plus bas.

Le deuxième cas de figure concerne la reconstitution des prénoms. Bien que celle-ci augmente la proportion de résultats valides, un prénom mal reconstitué risque de retourner du silence. C'est pourquoi une clé `rebuilt` a été ajoutée au dictionnaire qui contient les données extraites : cette clé permet d'indiquer qu'un nom a été reconstitué. Dans ce cas,

---

10. Dans ce qui suit, le terme « paramètre » se réfère aux clés du dictionnaire d'informations structurées, et non aux paramètres de l'API présentées ci-dessus (`srsearch...`).

une nouvelle recherche est lancée sans le prénom. Si aucun résultat n'est obtenu, alors c'est l'algorithme de recherches soustractives qui s'exécute.

Le troisième cas de figure concerne la présence de plusieurs paramètres de recherche disponibles – en langage technique, de plusieurs valeurs non-nulles dans le dictionnaire. Dans ce cas, c'est directement l'algorithme soustractif qui s'exécute. Celui-ci a été conçu à partir du constat que, d'un côté, les données contiennent inévitablement de bruit ; de l'autre, l'extraction de données risque elle aussi d'introduire du bruit. Cette logique soustractive cherche donc à supprimer le bruit qui a pu s'introduire. Afin d'éviter d'être trop « brut », l'algorithme procède en séparant les différents types de bruit possibles. D'abord, il arrive souvent qu'un chiffre soit mal reconnu lors de l'océrisation, et il peut donc y avoir des erreurs dans les dates. Deux dates peuvent être extraites (naissance et décès) ; il y a cependant peu de chances d'erreur sur les deux dates. Par conséquent, une des deux dates est supprimée et la recherche est lancée ; si aucun résultat n'est obtenu, c'est l'autre date qui est supprimée. La seconde source de bruit provient des prénoms : comme cela a été dit, au XIX<sup>ème</sup> s., les prénoms étrangers sont régulièrement traduits ; il arrive souvent que les prénoms composés ne soient pas dans le même ordre que sur *Wikidata*, ou qu'une personne soit nommée par son deuxième prénom – peut-être parce que les catalogues utilisent des noms d'usage alors que *Wikidata* s'en tient aux noms civils. Pour ces raisons, le prénom est retiré de la chaîne de caractères à rechercher avant de faire une requête sur l'API. Bien que les risques liés aux dates et aux prénoms ont été corrigés (autant que possible), il est possible que l'auteur.ice n'ait pas encore été aligné avec une entité *Wikidata*. À ce stade, tous les risques spécifiques ont été traités. L'algorithme relance alors des recherches en retirant à chaque fois un paramètre différent, jusqu'à ce qu'un identifiant ait été extrait de *Wikidata* ou que toutes les permutations possibles aient été faites.

Cette série de construction d'URL cherche à minimiser les redondances ; cependant, il est possible qu'une même chaîne de caractères puisse être extraite à différentes étapes. C'est pourquoi des fichiers de log sont créés en cours d'utilisation : ils servent à sauvegarder les chaînes de caractères qui ont déjà été recherchées et les résultats qui leur sont associées. Avant de lancer une requête sur l'API, l'algorithme vérifie donc si la chaîne a déjà été recherchée. Si oui, il récupère le résultat obtenu ; sinon, lance la recherche et sauvegarde le résultat retourné par *Wikidata*.

L'algorithme de recherche sur *Wikidata* est visiblement complexe, puisque son objectif est de minimiser le taux d'erreur possible sur des données variées, ce qui demande d'adopter des comportements différents en fonction des informations disponibles. L'algorithme ayant tendance à retirer de plus en plus de paramètres de recherches en cours d'exécution (les dates, les prénoms...), il est possible que la chaîne recherchée finisse par être très pauvre en informations ; cela permet de s'aligner avec une entité, mais celle-ci risque de ne pas correspondre au `tei:name` ou d'être très générique. Si l'entité recherchée

est « Napoléon Bonaparte », mais que suffisamment de paramètres ont été retirés, il est possible que l'algorithme ne recherche que « Bonaparte ». Le résultat obtenu est alors « Charles Lucien Bonaparte ». Il est donc utile d'avoir une mesure de la certitude (ou l'incertitude) de la validité d'un résultat, ne serait-ce que pour accélérer le processus de correction des identifiants (ce qui, sur plus de 82000 entrées, n'est pas une mince affaire). Le parti pris a été, tout simplement, d'établir un score de certitude en fonction du nombre de paramètres utilisés pour lancer une recherche. Ce score a été défini de façon expérimentale, en mesurant la proportion de résultats considérés comme certains et ainsi que le taux d'erreur au sein des résultats certains (c'est-à-dire, le nombre de résultats qui ont dépassé le seuil de certitude mais qui sont en fait erronés). Après plusieurs essais différents, le score de certitude  $c_r$  pour un résultat  $r$  a été déterminé. Il suit la formule ci dessous :

Étant donnés :

$c_r$  le score de certitude pour un résultat  $r$  ;

$q_r$  la chaîne de caractère recherchée pour obtenir le résultat  $r$  ;

$\sum param$  la somme des paramètres utilisés dans  $q_r$  ;

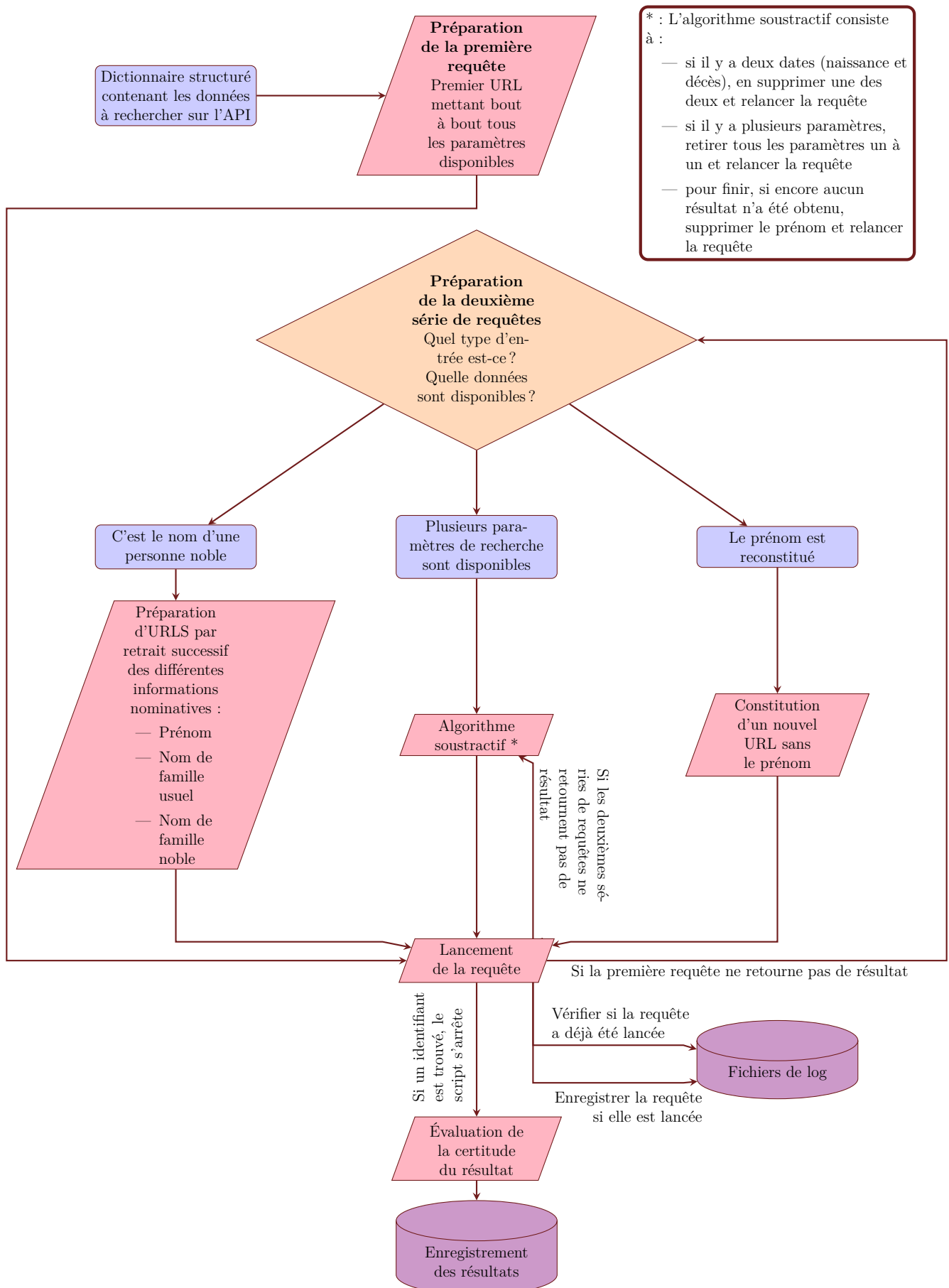
$d = 1$  si une date est dans  $q_r$  ; sinon,  $d = 0$  ;

$p = 1$  si le prénom utilisé dans  $q_r$  n'a pas été reconstitué ; sinon,  $p = 0$  ;

$$c_r = \sum param + d + p$$

Un résultat est considéré comme étant « certain » si  $c_r \geq 4$  ou si des dates sont présentes dans  $q_r$  (les chances qu'un alignement soit fait entre un `tei:name` et une entité homonyme ayant les mêmes dates de naissance ou de décès étant très faibles). Sur le jeu de données de test, 32% des entrées atteignent ce seuil de certitude. 23,6% de celles-ci comprennent en fait une erreur, soit 8,5% du jeu de données total (voir le tableau en annexes : 2). L'objectif de ce score de certitude n'est pas d'être parfaitement exact, mais d'accélérer la relecture des résultats : en acceptant un taux d'erreur total de 8,5%, il est possible de corriger que les entrées pour lesquelles  $c_r < 4$ . Seuls les deux tiers du jeu de données restent alors à corriger manuellement.

Cet algorithme cherche à répondre à plusieurs problèmes techniques : il travaille avec des données historiques contenant du bruit et un niveau très inégal d'information. Son fonctionnement vise à maximiser la probabilité d'obtenir un résultat valide, en faisant des requêtes d'abord très spécifiques puis de plus en plus génériques. Cependant, il crée lui-même plusieurs problèmes : son fonctionnement est assez complexe et demande de manipuler plusieurs lourds fichiers en même temps. Un ensemble d'optimisations ont donc été réalisées.

FIGURE 5.3 – Algorithme lançant des recherches en plein texte sur *Wikidata*

## 5.4 Gérer la montée en charge : optimisation et réduction du temps d'exécution

L'intégralité de l'algorithme (contenant extraction d'informations et interaction avec une API) a cessé de fonctionner sur mon ordinateur alors qu'à peine 5% des 820000 entrées avaient été traitées. Ce problème technique soulève des difficultés adjacentes à la recherche en humanités numériques : dans ce contexte, des ordinateurs très performants ne sont pas toujours disponibles, et les machines utilisées sont souvent des ordinateurs de bureau traditionnels. De telles situations, la recherche peut parfois faire face à des limitations matérielles. De telles limites sont – bien sûr – un problème ; en obligeant chercheur.euse.s à intégrer des contraintes matérielles dans leurs algorithmes, cette difficulté peut aider à réaliser des programmes plus efficaces, moins gourmands en ressources et en énergie, et donc aider à faire de la recherche plus durable. De telles considérations écologiques sont partagées par toute une communauté de chercheur.euse.s issu.e.s des humanités et défendant des approches « low-tech »<sup>11</sup> ; ces considérations écologiques sortent du champ des humanités et se retrouvent dans des travaux provenant de l'informatique « traditionnelle » et de l'apprentissage machine<sup>12</sup>. L'optimisation de programmes informatiques n'est pas seulement une question écologique : c'est également une manière d'écrire de meilleurs programmes, plus performants et efficaces.

Dans sa première version, le script d'alignement avec *Wikidata* avait un fonctionnement apparemment plus « simple », puisqu'il ne prenait pas du tout en compte les questions d'optimisation : pour chaque entrée de catalogue, des informations étaient extraites et les recherches faites sur *Wikidata*. Ce fonctionnement avait deux défauts. D'abord, il demandait que tout le script soit exécuté en une fois – si il y a une erreur imprévue, il est nécessaire de recommencer à nouveau. Ensuite, une interaction avec *Wikidata* était nécessaire pour chaque entrée, peu importe si le nom contenu dans le `tei:name` avait déjà été recherché auparavant (alors qu'une même personne est souvent l'auteur.ice de plusieurs manuscrits recensés dans les catalogues) ; l'interaction avec *Wikidata* prenant plusieurs secondes, cela peut grandement ralentir le processus d'alignement avec *Wikidata*. C'est à partir de ces deux problèmes que le processus d'optimisation a commencé.

Le premier problème – l'impossibilité de vérifier quelles entrées de catalogue ont déjà

---

11. Le *Minimal computing* (« Informatique minimale » en français) est un groupe de réflexion créé au sein de *GO ::DH – Global Outlook Digital Humanities*. Il s'intéresse aux manières de rendre l'informatique minimale, dans différents sens du terme : à la fois plus accessible, mais également plus écologique. Plus généralement, le courant du *minimal computing* développe une approche réflexive sur les façon de faire de la recherche durable et libre (au sens de l'*open source*) ; ce groupe s'intéresse également au développement des humanités numériques dans un contexte non-occidental. Jentery Sayers, *Minimal definitions*, Minimal Computing - a working group of GO ::DH, 2 oct. 2016, URL : <http://go-dh.github.io/mincomp/thoughts/2016/10/02/minimal-definitions/> (visité le 10/06/2022).

12. Emma Strubell, Ananya Ganesh et Andrew McCallum, « Energy and Policy Considerations for Deep Learning in NLP », *CoRR*, abs/1906.02243 (2019), DOI : <https://doi.org/10.48550/arXiv.1906.02243>.

été traitées – a d’abord été corrigé en consultant, à chaque fois qu’un nouvel identifiant était recherché, le fichier de sortie. Celui-ci est un tableur qui contient l’identifiant des entrées, le `tei:name`, le `tei:trait`, l’identifiant récupéré sur *Wikidata* et un extrait de la description associée à cet identifiant. Si cette solution permet de ne pas recommencer à zéro à chaque fois que le script est lancé, elle demande de charger en mémoire un très large fichier `csv` (12,7 mégaoctets). Cette vérification particulièrement gourmande a été remplacée par la création d’un fichier de log annexe en cours d’exécution. Celui-ci contient uniquement les identifiants des différentes entrées de catalogue et mesure seulement 1,3 mégaoctets. À chaque fois qu’une nouvelle entrée est traitée, ce fichier est ouvert, et le script vérifie que le nouvel identifiant ne s’y trouve pas. Si l’identifiant ne s’y trouve pas, alors le `tei:name` et le `tei:trait` de cette entrée sont traités et alignés avec un identifiant *Wikidata*. À la fin de ce processus, l’identifiant `TEI` de l’entrée est ajouté au fichier de log. Cette simple modification permet donc de vérifier si une entrée de catalogue a déjà été traitée à l’aide d’un fichier dix fois moins lourd que dans la tentative précédente.

Ce processus permet de vérifier quelles entrées ont déjà été traitées, mais il ne permet pas de savoir si une personne se retrouve d’une entrée sur l’autre. Si une telle récurrence est identifiée, il serait possible de récupérer les résultats déjà disponibles, plutôt que de réaliser à nouveau l’alignement avec *Wikidata*. Non seulement cela pourrait rendre le script plus rapide – le lancement des requêtes pour une entrée prenant environ 2 à 5 secondes –, mais en plus, cela éviterait de surcharger l’API de requêtes, ce qui, dans un cas extrême, peut faire saturer une application en ligne et perturber son fonctionnement. La difficulté consiste alors à désambiguïser les informations présentes dans le `tei:name` (c’est-à-dire, à repérer toutes les occurrences d’une même personne dans les `tei:name`). Cette désambiguïsation est compliquée (l’un des intérêts de la reconnaissance d’entités nommées est justement d’identifier les multiples occurrences d’une même personne). Plutôt que de trouver un moyen de désambiguïser réellement les entrées, une solution plus simple a été mise en place : chaque chaîne de caractère recherchée est enregistrée dans un fichier de log, avec l’identifiant *Wikidata* récupéré par recherche. L’intégralité du processus d’extraction d’information des `tei:name` et `tei:trait` est donc menée sur toutes les entrées ; toutes les chaînes de caractères à rechercher sont également construites ; cependant, une recherche n’est faite sur l’API que si elle n’a jamais été faite auparavant. Du fait de la taille du jeu de données, un problème apparaît cependant : plusieurs chaînes de caractères étant construites pour chaque entrée de catalogue, l’utilisation d’un seul fichier de log pour contenir toutes ces recherches amènerait à la création d’un immense fichier (plusieurs centaines de milliers de lignes). Il serait alors nécessaire de le parcourir intégralement à chaque fois qu’une requête est lancée ; le script optimisé risquerait alors d’être moins efficace qu’avant son optimisation. Le choix a donc été fait de créer automatiquement et d’utiliser plusieurs fichiers de log : les chaînes de caractères y sont enregistrées en fonction du premier caractère qu’elles contiennent, avec un fichier de log par lettre. Ainsi, avant de



lancer une recherche, le fichier à parcourir est bien plus court (2060 entrées pour la lettre « e »).

Avec toutes ces optimisations, le script a pu fonctionner en utilisant au maximum 3 gigaoctets, soit l'équivalent de naviguer sur internet en ayant un logiciel de traitement de texte ouvert (sur ma machine). Cependant, l'utilisation fichiers de log faite ici – qui a elle même été optimisée – demande de faire un arbitrage. En n'enregistrant pas de données et en lançant des requêtes à chaque fois, c'est la charge processeur qui augmente, ainsi que la consommation de courant due à internet. En enregistrant les requêtes déjà lancées, ces charges là sont diminuées, mais l'utilisation de mémoire vive augmente (il faut charger en mémoire et lire de larges fichiers). De ces deux deux différentes approches (augmenter la charge processeur ou augmenter la mémoire vive) il est nécessaire de déterminer laquelle est la plus efficace. C'est un des objectifs des tests menés sur l'algorithme d'alignement avec des identifiants *Wikidata*, comme nous allons le voir maintenant.

## 5.5 Évaluer la résolution d'entités : performance, qualité des données extraites de *Wikidata* et comparaison avec d'autres projets

Le processus d'alignement avec des entités *Wikidata* n'est pas parfait, et implique nécessairement un taux d'erreur. Au vu de la taille du corpus, la correction des résultats est un processus particulièrement chronophage. Pour pouvoir utiliser les données produites avant cette correction, il est nécessaire de trouver une manière d'évaluer la qualité des résultats. Des tests ont été réalisés pour faire cette évaluation. Ils ont également été utilisés durant tout le développement de ce processus, afin de comparer les différents algorithmes d'extraction et de requêtes entre eux. Les tests ont également permis de mesurer l'impact de chaque paramètre de recherche dans l'obtention d'un résultat valide. Enfin, l'évaluation des algorithmes permet de comparer sa version optimisée et non-optimisée.

Comme cela a déjà été dit, les tests ont été lancés sur un jeu de données représentatif contenant 200 `tei:name` et `tei:trait`, choisis dans différents catalogues ; ce jeu de test représentant également les différents types de noms possibles (de personnes nobles, de lieux...); enfin, la proportion de `tei:trait` dans le jeu de données complet et dans celui de test est la même. Deux catégories de tests ont été réalisés. Le premier, décrit plus haut, concerne l'impact de chaque paramètre dans l'obtention d'une entité correcte sur *Wikidata* (1). Le second sert à mieux identifier les performances de l'algorithme final (2).

Ce dernier test réalise l'intégralité du processus d'extraction d'informations et d'alignement avec *Wikidata* sur l'ensemble du jeu de test. Ce processus est réalisé trois fois sans utiliser des fichiers de log et trois fois en les utilisant (c'est-à-dire, en inscrivant les recherches déjà faites dans un fichier et en vérifiant avant chaque recherche si une chaîne

de caractère a déjà été recherchée). Faire plusieurs fois de suite le même processus permet de les variations de performance induites par l'utilisation de fichiers de log. La performance d'un script, en terme d'utilisation des processeurs et de la mémoire vive, est assez difficile à quantifier ; elle peut également varier d'une machine à une autre ; enfin, il n'est pas possible de mesurer la charge supplémentaire infligée à un serveur distant lors de l'utilisation répétée de sources de données en ligne. Le seul critère retenu a donc été le temps d'exécution. Le script étant lancé trois fois en utilisant des fichiers de log et trois fois sans, les deux temps d'exécution  $\bar{t}$  correspondent au temps moyen nécessaire à traiter les 200 entrées du jeu de test. Soit  $t_i$  le temps pris pour traiter une fois le jeu de test, la moyenne  $\bar{t}$  correspond à :

$$\bar{t} = \frac{1}{3} \sum_{i=1}^3 t_i$$

La difficulté pour mesurer les temps d'exécution  $t_i$  en utilisant des fichiers de log est que, en dehors des tests, ces fichiers sont de plus en plus volumineux ; le temps utilisé pour les parcourir augmente en conséquence. Pour que ces  $t_i$  soient représentatifs, de faux fichiers de log sont créés avant le lancement des tests : ceux-ci comprennent des fausses données qui correspondent environ aux résultats obtenus pour 30000 recherches lancées sur l'API (soit un peu moins que les 82000 entrées du jeu complet, en prenant en compte que certaines recherches ne retournent pas de résultat, ce qui allège d'autant le fichier à parcourir). Le temps d'exécution  $\bar{t}$  en utilisant des fichiers de log est de 88,3 secondes ; sans ces fichiers,  $\bar{t}$  est de 92,5 secondes. Si la différence entre les deux résultats n'est pas immense, l'utilisation de fichiers de log permet également la charge processeur et l'interaction avec *Wikidata*, ce permet d'autres améliorations au delà du seul temps de performance.

À part cette seule évaluation de performance, les tests sur l'algorithme final (2) servent à mesurer la qualité des résultats obtenus. La mesure de performance pour un processus d'alignement est le score F1, qui prend à la fois en compte la précision (le nombre de vrais positifs obtenus parmi tous les résultats obtenus) et le rappel (le nombre de vrais positifs obtenus parmi tous les vrais positifs). Celui-ci est de 67,4%, ce qui marque une nette amélioration (près de 20 points) de l'utilisation seule du prénom et du nom de famille usuel extrait du `tei:name`.

Ce score est également intéressant en comparaison avec d'autres projets de NEL, sous-domaine du traitement automatisé du langage qui vise à associer une entité avec une base de connaissances externes. Aicha Soudani et al. ont, en 2018, présenté les résultats d'un projet de NEL de lieux sur un corpus relativement proche du notre : des textes littéraires du XIX<sup>ème</sup> s. encodés en XML-TEI.<sup>13</sup> Les méthodes utilisées par ce projet diffèrent

---

13. Aicha Soudani, Yosra Meherzi, Asma Bouhafs, F. Frontini, C. Brando, Yoann Dupont et Frédérique Mélanie-Becquet, « Adaptation et évaluation de systèmes de reconnaissance et de résolution des entités nommées pour le cas de textes littéraires français du 19ème siècle », dans *Atelier Humanités Nu-*

largement de celles présentées ici : les auteur.ice.s s'appuient fortement sur l'apprentissage machine, en identifiant d'abord des entités dans le texte à l'aide de SEM avant de lier les entités à différentes bases de connaissance en utilisant REDEN. Leur projet présente donc une complexité supplémentaire, puisqu'il demande de repérer les entités, ce qui n'est pas le cas ici (elles se trouvent toutes dans des balises `tei:name`). Cependant, c'est cette deuxième étape qui nous intéresse ; les entités ont déjà été identifiées, ce qui nivelle cette différence entre les projets. Trois bases de données en ligne ont été utilisées : *DBPedia*, *DataBnF* et *Wikidata*. En comparaison avec ce projet, il est intéressant de remarquer que le processus d'alignement avec *Wikidata* développé par *MSS / Katabase* obtient un score F1 très satisfaisant (67,4%). Celui-ci dépasse les scores obtenus par Soudani et al. dans leur alignement avec *Wikidata* (61,1%) et *DBPedia* (53,6%) (voir le détail des résultats en annexe : 3). Les résultats comparativement faibles de Soudani et al. s'expliquent en partie par le fait que REDEN recherche des candidats pour le liage d'entités en établissant une correspondance parfaite des caractères (*exact string match*)<sup>14 15</sup>. Les scores comparativement positifs du projet *MSS / Katabase* sont peut être également rendus possibles par la quantités d'informations disponibles sur chaque personne dans le corpus, récupérée grâce à un travail d'extraction des données sur l'ensemble du `tei:trait` et du `tei:name`.

## 5.6 En conclusion : retour sur l'extraction d'informations des catalogues et sur l'algorithme d'alignement avec *Wikidata*

L'algorithme liant les noms d'auteur.ice dans les catalogues avec des entités *Wikidata* propose une réponse à plusieurs difficultés techniques propres à la reconnaissance d'entités par la détection de motifs sur un corpus de données varié. Parmi toutes ces difficultés possibles, les trois les plus importantes sont le bruit dans les données d'origine, le bruit produit par l'extraction de données et enfin la difficulté à réaliser un alignement entre des documents historiques et des bases de données contemporaines. La solution proposée s'appuie sur un fonctionnement modulaire, qui s'adapte aux données produites et prend en compte la possibilité que du bruit se soit introduit à une étape ou à une autre. De différence de résultats obtenue avec Soudani et al., il faut retenir que la technique utilisée n'est pas toujours garante de la qualité des résultats obtenus. Plutôt que de chercher à établir des correspondances exactes, l'approche décrite dans ce mémoire, basée sur une

---

*mériques Spatialisées (HumaNS'2018)*, Montpellier, 2018, URL : <https://hal.archives-ouvertes.fr/hal-01925816> (visité le 10/08/2022).

14. *Ibid.*, p. 4.

15. Il est également à remarquer que Soudani et al. disposent d'un score d'« exactitude totale du liage » (*Overall linking accuracy*) qui est nettement plus élevé (entre 70 et 85%) ; la méthode de calcul de ce score, qui mesure la performance de l'ensemble du processus de NEL, n'est cependant pas décrite dans l'article. Il est donc impossible d'y comparer les résultats de *MSS / Katabase*.

détection de motifs *low-tech* a permis de développer une méthode très précisément adaptée à la fois au corpus et à *Wikidata*. Ce qui est potentiellement perdu par la détection de motifs est donc corrigé par l'algorithme menant des requêtes sur *Wikidata*.

L'alignement avec des entités *Wikidata* est la meilleure manière de corriger le bruit dans les données du projet : en liant les différents `tei:name` à des identifiants *Wikidata*, il devient possible d'extraire des données structurées d'une source externe, pour construire une base de connaissance plus fiable : même si un corpus semi-structuré facilite le travail d'extraction d'informations, l'accès aux informations reste relativement complexe ; il existe de plus des biais dans les informations contenue par le corpus (ce qui est mis en avant dans les `tei:trait` est avant tout ce qui est susceptible d'intéresser le public de l'époque). En se servant de ce liage d'entités pour constituer une base de données, il est possible de chercher à corriger ce biais dans les données, tout en rassemblant assez d'informations pour pouvoir faire une analyse des facteurs déterminant le prix d'un manuscrit sur le marché du XIX<sup>ème</sup> s.. C'est la manière dont est constitué le jeu de données issu de *Wikidata* qui est expliqué ci-dessous.

# Chapitre 6

## Après l’alignement, l’enrichissement : utiliser SPARQL pour produire des données structurées

Le liage des noms contenus dans le `tei:name` avec *Wikidata* n’est qu’une étape, et non la fin en soit du processus – bien que ce soit la partie la plus difficile. Les identifiants ainsi récupérés servent de « porte d’entrée » aux données disponibles sur la base de connaissance : une fois que ces identifiants sont connus, il est possible de récupérer automatiquement via SPARQL des données sur les auteur.ice.s et autres entités présentes dans les catalogues. Obtenir ces informations permet avant tout de mieux situer l’auteur.ice d’un manuscrit, afin mieux le statut d’une personne et de comprendre quelles informations biographiques influencent le prix d’un manuscrit. Mais le liage d’entités nommées peut avoir d’autres fins et permettre d’enrichir le corpus de catalogues de plusieurs manières, comme nous le verrons.

### 6.1 Comprendre les particularités des modèles sémantiques de données

Un bref retour sur la manière dont s’organisent les données contenues dans *Wikidata* permet de mieux comprendre le problème que peut poser la diversité du corpus pour la constitution d’une base de connaissances propre au projet. *Wikidata*, comme de nombreuses autres bases en lignes (*DataBnF*, *DBPedia*...), permet d’accéder à des informations stockées dans une base de données en graphe. Ce type de base de données a un fonctionnement très particulier qui influence la constitution d’une base de connaissance à partir des identifiants *Wikidata* récupérés à l’étape précédente. Dans une base de données

en graphe, les données sont encodées en XML-RDF. Ce format – dit sémantique – contient des données liées sous forme de « triplets » sujet–prédicat–objet, où :

- le sujet est la ressource principale.
- le prédicat est une propriété du sujet, qui caractérise une relation avec une autre ressource, l'objet.
- l'objet est une ressource secondaire : c'est la valeur d'un prédicat.

Le principe des triplets RDF est mieux exprimé sous forme graphique (6.1) :

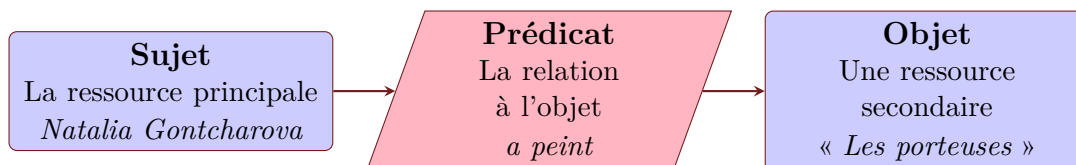


FIGURE 6.1 – Exemple de relation sujet – prédicat – objet

Trois particularités supplémentaires définissent les formats sémantiques :

- Toutes les « ressources » peuvent être tour à tour sujet ou objet. L'exemple du dessus, par exemple, aurait pu être réécrit sous la forme : « Les porteuses » a été peint par Natalia Gontcharova. Dans ce cas, Natalia Gontcharova est l'objet et *Les porteuses* est le sujet. Par conséquent, une base de données RDF est une base de données en graphe ; elle peut être représentée sous la forme d'un réseau de ressources qui entretiennent des relations bilatérales entre elles. Il n'y a pas de hiérarchie entre les informations, contrairement à une base de données XML classique.
- L'ensemble des ressources et des prédicats d'une base de données en graphe sont définis et disposent d'un identifiant unique (certains de ces identifiants ont été récupérés grâce au liage d'entités nommées).
- Si le langage SPARQL offre une syntaxe commune, chaque base de connaissance peut utiliser des vocabulaires particuliers qui sont organisés en « ontologies ». Une ontologie correspond à la définition d'un ensemble de catégories, de propriétés et de relations qui unissent des données et des concepts ; cet ensemble est complété par une modélisation (souvent sous forme graphique), qui indique la relation entre les différents termes de l'ontologie<sup>1</sup>. Ceux-ci sont souvent liés de façon hiérarchique (plusieurs termes spécifiques pouvant être dérivés d'un terme générique). Au sein d'une ontologie, chaque terme a lui-même un identifiant unique, ce qui garantit une implémentation uniforme pour l'ontologie. Une même base de connaissance peut utiliser plusieurs ontologies. Celles-ci sont définies à l'aide d'espaces de noms<sup>2</sup>,

1. *Ontology (information science)*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://en.wikipedia.org/wiki/Ontology\\_\(information\\_science\)](https://en.wikipedia.org/wiki/Ontology_(information_science)) (visité le 01/02/2022).

2. *Namespace*, Wikipedia. L'encyclopédie libre, 2022, URL : <https://en.wikipedia.org/wiki/Namespace> (visité le 05/05/2022).

c'est-à-dire par des identifiants qui permettent de différencier les ontologies. Les prédicats, plus particulièrement, sont définis selon une ontologie particulière.

SPARQL a l'avantage de permettre de récupérer des données propres sur des bases en ligne et offre une syntaxe unique partout où il est implémenté. Cependant, le troisième point complexifie son utilisation, ainsi que la définition des données à récupérer : les prédicats sont décrits avec une grande précision et selon des vocabulaires spécifiques ; par conséquent, une information analogue peut être représentée par différents prédicats selon le type de donnée qui est requêtée (une personne, une sculpture...). Dans l'ontologie *Wikidata*, la création d'un texte et la création d'une peinture ne correspondent pas au même prédicat. Pour que les données soient utilisables, il faut être très spécifique quant aux informations recherchées.

## 6.2 Quelles données rechercher via SPARQL ?

La question des données qui doivent être récupérées, et donc de la base de connaissance à constituer pour le projet *MSS / Katabase* n'est pas anodine. D'abord, l'utilisation de sources externes peut être utilisée pour corriger des biais dans les données originelles (où ce sont les auteurs, éditeurs et éditrices de catalogues qui décident des informations à intégrer). Ensuite, si l'objectif premier est de récupérer des données pour mener une étude économétrique, le liage avec *Wikidata* peut également permettre d'autres enrichissements et fonctionnalités. Enfin, cette étape représente une difficulté technique : le corpus de catalogues et les entités récupérées sont très diverses, tandis que SPARQL utilise au contraire des vocabulaires très spécifiques. Il faut donc construire des requêtes très étendues, afin d'obtenir des résultats pour tous les types de données. Les 18899 entités avec lesquelles les entrées de manuscrits ont été alignées peuvent se classer en de nombreuses catégories. Sur *Wikidata*, une entité est une « instance » d'une classe plus large. C'est en fonction de ces classes qu'il faut construire les requêtes : les propriétés recherchées pour chaque entité peuvent être spécifiques à une classe, mais pas à l'entité elle-même. En suivant la classification de *Wikidata*, les entités présentes dans le corpus appartiennent aux catégories suivantes :

- personnes humaines ; cette catégorie est la plus fréquente (12090 occurrences)
- noms de familles (3180 entités)
- communes françaises (586 occurrences)
- peintures et sculptures (respectivement 520 et 236 entités)

Cette variété – présentée sous forme graphique en annexes (1) – peut s'expliquer en partie par le taux d'erreur dans l'alignement avec *Wikidata* : il est par exemple probable que les peintures et sculptures soient sur-représentées parmi les entités *Wikidata* liées

aux catalogues. Cependant, toutes ces « erreurs » ne correspondent pas forcément à des résultats qui ne sont pas pertinents. Par exemple, l'algorithme peut aligner un.e écrivain.e avec un de ses ouvrages, ou une personne avec son portrait. Des résultats erronés peuvent toujours garder une forme de pertinence. Il est d'autant plus important de construire des requêtes SPARQL qui se concentrent pas uniquement sur des personnes. Cependant, il n'est pas possible de s'adapter à l'intégralité de la diversité du corpus. Le choix a donc été fait de se concentrer sur les catégories les plus pertinentes : les personnes, les familles, et les œuvres artistiques et littéraires. Non seulement ces catégories contiennent la grande majorité du corpus (16026 entités, soit 84,79%), mais ces catégories sont les plus à même de contenir des entités pertinentes. Il a été choisi de ne pas faire de requête spécifique sur les lieux, puisque *Wikidata* présente peu d'informations pour les entités de la catégorie « communes françaises ».

Un nombre assez conséquent de données ont donc été requêtées avec SPARQL, du fait des spécificités des bases de données en graphes, de la variété des entités *Wikidata* auxquelles les manuscrits sont liées, et enfin du fait de la variété du corpus lui même. Ces informations récupérées correspondent aux différentes catégories de *Wikidata*.

- Pour les personnes et les familles, les informations suivantes sont récupérées sur *Wikidata* :
  - Le genre de la personne ;
  - Sa nationalité, afin de voir si l'origine d'une personne influence le prix d'un manuscrit ;
  - Les langues parlées par une personne ; là encore, l'objectif est d'étudier l'impact de l'origine d'un.e auteur.ice sur un prix.
  - Les date de naissance et de décès, afin de placer un manuscrit dans une époque et de voir comment son ancienneté ou sa contemporanéité en influencent le prix.
  - Le lieu où une personne est née, où elle a vécu et où elle est morte, pour des raisons analogues.
  - La manière dont la personne est morte. Si cette information peut sembler anecdotique à un public contemporain, les catalogues de ventes sont marqués par un goût du sensationnel, et la manière dont une personne est morte est souvent mentionnée, notamment en cas d'exécutions.
  - La religion d'une personne : il peut être intéressant d'étudier si, et comment, ce critère influence l'évolution d'un prix.
  - Les titres de noblesse d'une personne.
  - L'éducation qu'a reçu une personne, afin de mieux situer ses occupations et d'analyser l'impact du niveau et du type d'éducation sur le prix.



- L’occupation d’une personne, et les fonctions précises qu’elle a occupées : là encore, il est intéressant de situer l’impact de la carrière sur le prix et de voir quelles occupations sont corrélées avec des prix élevés sur le marché des manuscrits.
- Les prix et distinctions reçus par une personne. À l’aide de ce critère, il est alors possible de chercher à répondre à cette question : la célébrité d’une personne de son vivant impacte-elle le prix de ses manuscrits ?
- Les organisations et institutions dont la personne est membre (Académie française, Franc-maçonnerie...)
- Le nombre d’œuvres écrites ou réalisées par une personne. Là encore, c’est une tentative de mesurer l’impact ou la célébrité d’une personne : les manuscrits de quelqu’un ayant beaucoup écrit sont-ils plus chers que les manuscrits d’une personne ayant peu écrit ?
- Le nombre de conflits auxquels une personne a participé. Ce critère de recherche permet de quantifier l’importance d’un personnage militaire.
- Des images, telles que le portrait et la signature.
- Pour les créations littéraires, ce sont des informations bibliographiques qui sont avant tout récupérées ; pour les autres œuvres d’art, des informations analogues sur le contexte de création sont retenues.
  - Le titre de l’œuvre.
  - Son auteur.ice, pour étudier si certain.e.s auteur.ice.s sont susceptibles d’influencer le prix d’un manuscrit.
  - La date de création de l’œuvre, afin de savoir si l’époque d’origine influence le prix. Pour les livres, la date de publication est également récupérée.
  - La requête récupère aussi la maison d’édition d’un livre.
  - Les dimensions et matériaux d’une œuvre d’art sont également d’intérêt.
  - Enfin, le genre et le mouvement dans lequel s’inscrit une œuvre sont d’intérêt : ces informations pourraient permettre de voir si une hiérarchie des goûts influence le prix d’un manuscrit.
- Pour finir, afin de pouvoir éventuellement enrichir nos données avec d’autres sources externes à *Wikidata*, des identifiants uniques ont été récupérés afin de donner accès à d’autres bases de données en ligne : les identifiants VIAF (Fichier d’autorité international virtuel), ISNI (International Standard Name Identifier), de la Bibliothèque nationale de France, de la Bibliothèque du Congrès américain, ainsi que les identifiants IDRef. Certaines institutions, comme la BnF, rendent leurs données accessibles via SPARQL ; la récupération de ces identifiants faciliterait grandement les enrichissements ultérieurs depuis d’autres sources de données.

Comme on l'a dit, l'objectif principal de l'alignement avec *Wikidata* est de produire des données pour calculer des régressions linéaires, ce qui permettrait d'étudier les déterminants du prix d'un manuscrit sur le marché du XIX<sup>ème</sup> s.. Cette récupération d'informations en masse ouvre d'autres possibilités. Entre autres, de nombreuses données géographiques ont été récupérées (lieu de naissance, de décès, d'inhumation et de résidence). Il est ensuite possible de récupérer les coordonnées de ces lieux, afin de construire une cartographie des auteur.ice.s dont les manuscrits circulent sur le marché du XIX<sup>ème</sup> s. parisien. C'est d'ailleurs souvent à des visées cartographiques et de géoréférencement que sont menées des campagnes de reconnaissance d'entités nommées ; dans les dernières années en France, de nombreuses études ont été menées pour développer des cartographies à partir de textes littéraires français encodés en TEI (Soudani et al.<sup>3, 4</sup>). En développant une approche cartographique, le projet *MSS / Katabase* pourrait apporter de nouvelles possibilités pour de telles études : le corpus de catalogues étant une source secondaire sur l'histoire littéraire française, une approche géographique permettrait d'étudier les origines géographiques des auteur.ice.s, plutôt que d'analyser les lieux représentés dans leurs œuvres. Cette possibilité n'est pas anodine, puisqu'elle permettrait de mettre en relation la « parisianité » avec la construction du canon littéraire à Paris. Il serait également possible d'étudier la circulation des productions culturelles, et leur rayon d'influence. En croisant les données géoréférencées avec des informations chronologiques (dates de naissance et de mort...), ces questions peuvent également être étudiées de façon historique : comment l'influence de l'origine géographique sur la réception d'une œuvre évolue au fil des siècles ? Répondre à ces questions n'a pas été possible dans le cadre de mon stage ; cependant, grâce à l'enrichissement de données via SPARQL, il de telles études deviennent possibles, et les données pour mener ces analyses sont au moins en partie déjà disponibles. Produire des informations normalisées et exploitables pour la recherche implique donc de produire des données qui puissent être réutilisées avec d'autres problématiques de recherches.

### 6.3 Présentation générale

L'algorithme de récupération des informations sur *Wikidata* est considérablement plus simple que ce qui a été présenté lors de l'extraction d'informations des catalogues et de l'alignement avec des entités *Wikidata* : lors de ces étapes, les principales difficultés résultaient du bruit dans les données et de la nature « semi-structurée » des entrées de catalogues, qui demandait de s'adapter à de nombreux cas de figure. Une fois que des

---

3. Id., « Adaptation et évaluation de systèmes de reconnaissance et de résolution des entités nommées pour le cas de textes littéraires français du 19<sup>ème</sup> siècle »..., p.4-5.

4. Francesca Frontini, C. Brando, Marine Riguet, Clémence Jacquot et Vincent Jolivet, « Annotation of Toponyms in TEI Digital Literary Editions and Linking to the Web of Data », *Matlit Revista do Programa de Doutorado em Materialidades da Literatura*, 4-2 (11 juil. 2016), p. 49-75, DOI : 10.14195/2182-8830\_4-2\_3, p. 63-66.

identifiants *Wikidata* ont été extraits, le processus est bien plus simple – comme cela apparaît dans la figure 6.2. Un identifiant étant unique et servant de clé d’entrée à une base de données en graphe très structurée, un comportement uniforme peut être adopté pour récupérer des données issues de *Wikidata*. Il ne s’agit plus que, pour chaque identifiant *Wikidata* récupéré, de lancer des requêtes SPARQL et d’en stocker le résultat dans un fichier. Le comportement de l’algorithme est donc toujours le même.

L’algorithme commence par vérifier si un identifiant a déjà été traité dans un fichier de log ; l’utilisation de ces fichiers évite d’avoir à recommencer la récupération d’informations de *Wikidata* à chaque interruption du script, ce qui est essentiel puisque cette étape dure plus de dix heures. Si l’identifiant n’a pas été traité, alors plusieurs requêtes SPARQL sont lancées sur *Wikidata*. Les résultats sont retournés en JSON ou XML dans des formats définis dans une spécification du W3C<sup>5</sup>. Ceux-ci ont un but descriptif (ils décrivent la requête et les données retournées) ; les documents SPARQL sont donc très complets et peu malléables ; c’est pourquoi ils sont transformés en une base de connaissance au format JSON. Celle-ci, comme nous le verrons, ne contient que les données obtenues dans un format plus malléable.

Les requêtes SPARQL sont faites sur un serveur distant via le protocole HTTP et demandent parfois au serveur de traiter de très grandes quantités de données. Des erreurs peuvent donc avoir lieu. Pour permettre au script de fonctionner malgré ces erreurs, un système de gestion des erreurs a donc été mis en place. Par défaut, il est demandé au serveur de *Wikidata* de fournir les résultats de la requête en JSON, format très léger et malléable. Deux erreurs peuvent alors arriver : soit le JSON retourné par le serveur est mal formé (et ne peut donc être traité), soit la requête met plus excède la durée maximale autorisée (qui est d’une minute). Dans ce cas, la requête est lancée une seconde fois au serveur, mais cette fois-ci le format de réponse demandé est le XML ; cela permet d’éviter que le document soit mal formé ; il est également possible que la requête s’exécute dans le temps imparti la seconde fois. Si une erreur a encore lieu à la seconde requête, alors une entrée vide est associée à l’identifiant *Wikidata* dans la base de connaissance, pour signifier qu’aucune réponse n’a pue être obtenue.

---

5. Dave Beckett, Jeen Broekstra et Sandro Hawke, *SPARQL Query Results XML Format (Second Edition)*. W3C Recommendation 21 March 2013, W3C, 21 mars 2013, URL : <https://www.w3.org/TR/2013/REC-rdf-sparql-XMLres-20130321/> (visité le 10/05/2022).

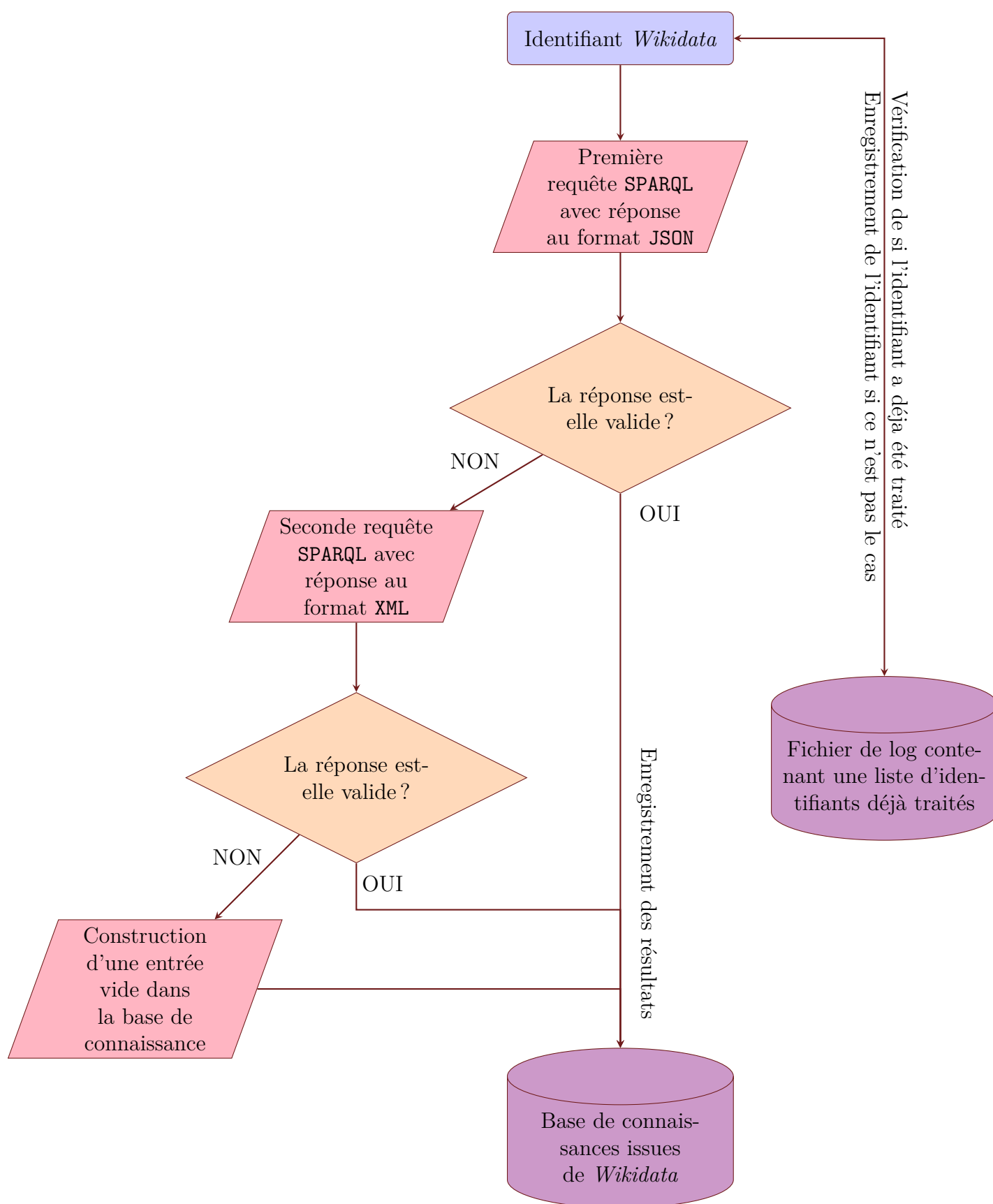


FIGURE 6.2 – L'algorithme de constitution d'une base de connaissances issues de *Wikidata* à l'aide de SPARQL

## 6.4 Développer un comportement uniforme pour produire des données exploitables à partir un corpus hétérogène

### 6.4.1 Les spécificités techniques de SPARQL : adapter les résultats à ses besoins

D'un point de vue technique, SPARQL suit la logique des « triplets » présentée ci-dessus. Le langage permet d'interagir avec une base de données en graphe au format RDF (récupérer des données, en envoyer, en mettre à jour ou en supprimer). Récupérer des données via SPARQL revient à demander au langage d'extraire de la base de données l'ensemble des valeurs qui correspondent à un triplet sujet-prédicat-objet. Par exemple, récupérer tous les tableaux peints par Natalia Gontcharova est possible à l'aide de la requête suivante (6.1) :

```

1 SELECT ?painting ?paintingLabel
2 WHERE {
3     ?painting wdt:P170 wd:Q232391 .
4
5     SERVICE wikibase:label {
6         bd:serviceParam wikibase:language "en" .
7     }
8 }
```

Code source 6.1 – Une requête simple : les identifiants et les noms de tableaux peints par Natalia Gontcharova

La requête ci-dessus consiste à récupérer tous les résultats possibles pour la variable `?painting` qui ont la relation P170 (« a pour créateur.ice ») avec l'entité Q232391 (l'identifiant de N. Gontcharova) : la requête récupère toutes les entités de *Wikidata* ayant l'artiste pour créatrice. La partie suivant `SERVICE...` permet d'associer à chaque `?painting` une variable `?paintingLabel` qui contient le titre de l'entité `?painting` en anglais. Chaque entité ou relation dispose d'un préfixe (`wdt:`, `wd:`, `wikibase:`) qui permet de lier une information à un espace de noms spécifique. De cette brève description du fonctionnement technique de SPARQL, il faut retenir que ce langage permet de récupérer d'une base de données toutes les informations qui correspondent à la proposition énoncée. Ces informations sont liées à des variables, ce qui permet d'accéder aux résultats renvoyés par SPARQL.

À l'origine, toutes les informations étaient demandées à SPARQL dans une seule requête. Celle-ci fonctionnait sans difficultés pour certaines entités ; dans d'autres cas cependant, elle ne s'exécutait pas dans le temps imparti. Plusieurs tentatives d'optimisation

ont été menées, mais sans succès sur d'autres. Ce cas de figure, comme cela est noté<sup>6</sup> par la documentation SPARQL de *Wikidata*, est systématique lorsqu'un très grand nombre d'informations sont récupérées pour une variable requêtée. Cela fait exploser la complexité du traitement de la requête par le serveur et cause un dépassement du temps maximal de requête autorisé. Le problème, c'est qu'il est difficile de prévoir les cas de figure où trop de résultats seront récupérés par SPARQL. Le choix a donc de continuer d'adopter un comportement uniforme. Plutôt que d'adapter les requêtes en fonction du type d'entités (et du risque qu'elle rende une requête trop difficile à traiter), plusieurs petites requêtes plus petites requêtes sont construites à la suite ; leurs résultats sont agglomérés afin qu'une entité corresponde à une série de résultats. De plus, les risques de dépassement du temps maximal de traitement autorisé par *Wikidata* ont été pris en compte : si cette durée est dépassée, la requête est relancée ; si aucune réponse n'est obtenue cette seconde fois, alors une série de résultats vides sont associés à l'entité problématique dans la base de connaissances.

La réponse retournée par *Wikidata* ne peut être utilisée en tant que telle et doit donc être retravaillée. En effet, en plus de contenir les données renvoyées par SPARQL, les formats de réponse définis par le W3C<sup>7</sup> décrivent la requête (en stockant un ensemble de variables) et les données retournées (en définissant systématiquement un type de données, et en associant à ce type une URI). Comme cela a été dit, une requête SPARQL associe une variable à une « question » ; celle-ci permet d'associer une ou plusieurs valeurs à cette variable. Une requête comprenant plusieurs variables peut donc retourner un nombre de valeurs différent pour chacune de ces variables. SPARQL ne peut alors pas déterminer de lien logique entre les valeurs obtenues pour ces variables. SPARQL fait donc le produit cartésien entre toutes les valeurs possibles pour chaque variable : il y a alors autant de « solutions à la requête »<sup>8</sup> qu'il y a de combinaisons possibles de résultats. Par exemple, si trois variables sont requêtées, que SPARQL retourne trois résultats pour chaque variable et qu'il n'existe aucun lien logique entre elles (la valeur d'une variable ne dépend pas de la valeur des autres), alors le nombre de « solutions à la requête » est de  $3 \times 3 \times 3 = 27$ . Des exemples de réponse aux formats JSON (13) et XML (14) sont disponibles en annexes. Dans les deux cas, la structure est la même :

- Un en-tête (**head**) contenant une liste des variables requêtées.
- Un corps de texte (**results**) contenant une liste de toutes les « solutions de requêtes possibles » (**bindings**), c'est-à-dire toutes les combinaisons possibles des valeurs obtenues pour les requêtées. Chacune de ces solutions contient une entrée par variable obtenue. Pour chaque variable, plusieurs informations sont retournées :

---

6. *Query optimization*, Wikidata, 2022, URL : [https://www.wikidata.org/wiki/Wikidata:SPARQL\\_query\\_service/query\\_optimization](https://www.wikidata.org/wiki/Wikidata:SPARQL_query_service/query_optimization) (visité le 06/06/2022).

7. Id., *SPARQL Query Results XML Format (Second Edition)*...

8. *Ibid.*, §2.3.1. Variable Binding Results.

- La valeur obtenue
- Le type de cette valeur, qui permet de classer la valeur en grands groupes : une URI, du texte libre...
- Si le résultat obtenu est en texte libre, alors la langue de ce texte est également fournie.

Cette variété d'informations documente précisément les données associées à chaque variable; cependant, cette structure imbriquée rend l'accès aux informations très difficile. Par conséquent, le format de réponse de **SPARQL** présente deux problèmes. D'abord, il contient des informations inutiles pour la base de connaissance qui doit être constituée : celle-ci n'a pas besoin du langage ou du type de valeurs récupérées associées à chaque variable. Ces informations étant définies dans la requête, leur valeur est déjà connue et peut être documentée dans alourdir la base de connaissance à construire. Ensuite, une réponse **SPARQL** peut potentiellement contenir énormément de redondances. Dans les faits, il n'est pas nécessaire de déterminer un lien entre les différentes variables requêtées, puisque celles-ci sont indépendantes. Ce qui importe, ce n'est pas le lien entre les distinctions obtenues par une personne et les positions qu'il ou elle a occupé (cette relation n'existe pas), mais la relation entre l'entité et les valeurs récupérées : ce qui est nécessaire, c'est de savoir que telle valeur obtenue entretient telle relation avec l'entité qui est requêtée (par exemple, de savoir que « Membre de l'Académie Française » est une distinction honorifique).

### 6.4.2 La base de connaissances à construire

L'objectif de cette base de connaissance est de stocker toutes les informations obtenues avec **SPARQL** pour toutes les entités, afin de faire les liens entre les auteur.ice.s de manuscrits et une source de données externe. Par conséquent, cette base de connaissance doit remplir les conditions suivantes :

- Permettre un accès facile aux données : cette base n'est pas censée être accessible au public, mais être facilement manipulable par une machine.
- Être construite dans un format léger et manipulable. Là encore, la base n'est pas censée être qualitative scientifiquement (contrairement à une édition textuelle en **TEI**), mais être utilisable rapidement par une machine. Le format utiliser doit donc favoriser le traitement des données, et non leur publication.
- Limiter au maximum les redondances d'informations : les requêtes contenant plus de 18000 identifiants, la base constituée sera probablement volumineuse; plus elle occupe de mémoire, plus il est coûteux (en temps et en charge processeur) d'y récupérer des données.

- Être malléable : le format et le schéma définis pour la base de données doit permettre de contenir, pour chaque variable, un nombre indéfini de valeurs.
- Être accessible depuis les identifiants des entités *Wikidata* : c'est eux qui permettent le lien avec les catalogues.

Trois formats ont été considérés pour la base de connaissance :

- Le tableur (type **TSV**) : connu et utilisé par tout.e chercheur.euse, il a l'avantage d'être très facilement compréhensible et consultable par un public non-spécialiste. Dans ce cas, chaque colonne correspondrait à une variable et chaque ligne à une entité, ce qui permet à la fois un accès rapide aux identifiants et aux résultats de requête.
- Une base de données document au format **XML**. Ce langage à balises a l'avantage d'être hiérarchisé et de permettre des arborescences très complexes. Il est de plus possible de s'aligner sur des standards, comme le **XML-TEI** ; de plus, il est possible de concevoir des schémas de validation pour un fichier **XML** ; ceux-ci garantissent qu'il n'y ait pas d'erreur dans la structure de données, ce qui peut être intéressant puisque la base de connaissance est construite de façon automatique, et qu'une erreur est donc potentiellement possible.
- Une base de données **JSON**. Ce format, également hiérarchisé, permet de construire des structures arborescentes aussi complexe que le **XML**. Il est par contre bien plus léger et facile d'utilisation que ce dernier ; contrairement aux langages à balises, le balisage sémantique est moins facile (il est moins évident d'associer une signification à chaque élément). De plus, il n'est à ma connaissance pas possible de construire des schémas de validation pour un **JSON**.

C'est le **JSON** qui correspond le mieux aux besoins du projet. Étant donné qu'il est impossible de prévoir le nombre de valeurs obtenues pour une variable à chaque requête, un format non-hiérarchisé est peu intéressant : dans un tableur, il n'est possible de stocker qu'une information par colonne/variable, à moins de créer des subdivisions au sein d'une colonne, ce qui rendrait le tableur plus complexe à manipuler. Au contraire, un format hiérarchisé permet d'associer chaque variable à une liste de résultats dont la longueur n'est pas définie par avance. Le **JSON** a également été privilégié au **XML** : ce dernier est considérablement plus complexe à manipuler, que ce soit pour y ajouter des informations que pour accéder à celles-ci. De plus, le langage à balises est bien plus lourd, et donc potentiellement plus long à manipuler. Étant donné que la base de connaissances est un format de « travail », et non de publication, un balisage sémantique n'est pas nécessairement intéressant.

La structure suivante a été définie pour la base de connaissance : à chaque identifiant *Wikidata* est lié les résultats obtenus via **SPARQL**. Ces résultats correspondent à toutes les variables requêtées ; à celles-ci sont associées une liste de résultats. Toutes les variables



requêtées sont toujours présentes. Si aucun résultat n'a été trouvé pour une variable, alors celle-ci est associée à une liste vide. Cela permet que toutes les entrées de la base aient toujours la même structure. Ainsi, la manipulation de la base est facilitée ; en contrepartie, celle-ci est plus lourde que si les seules variables présentes étaient celles pour lesquelles des résultats ont été trouvés. En d'autres termes, la base de connaissance correspond au schéma ci-dessous (6.2). Un exemple contenant deux entrées de la base figure en annexes (15)

```

1 {
2   "identifiant wikidata 1": {
3     "variable 1": ["résultat 1", "résultat 2", "résultat N"],
4     "variable 2": ["résultat 1", "résultat 2", "résultat N"],
5     "variable N": ["résultat 1", "résultat 2", "résultat N"],
6   },
7   "identifiant wikidata N": {
8     "variable 1": ["résultat 1", "résultat 2", "résultat N"],
9     "variable 2": ["résultat 1", "résultat 2", "résultat N"],
10    "variable N": ["résultat 1", "résultat 2", "résultat N"],
11  },
12  # autres entités...
13 }
```

Code source 6.2 – Structure de la base de connaissance construite à partir des résultats obtenus de SPARQL

### 6.4.3 Des réponses SPARQL à la base de connaissance

Pour produire des données utilisables dans un contexte de recherche, il est donc nécessaire de retraiter le format de réponse de SPARQL pour s'adapter au modèle de données défini ci-dessus. Cependant, le contenu d'une réponse SPARQL varie selon la requête lancée, selon l'identifiant *Wikidata* requêté et enfin selon les informations retournées par cette base de connaissance. En effet, une variable requêtée n'apparaît dans le corps des résultats que si des valeurs correspondant à cette variable se trouvent dans *Wikidata*. Il faut donc trouver une méthode uniforme pour traiter toutes les réponses de SPARQL, peu importe la requête lancée ou les résultats obtenus. C'est dans l'en-tête des réponses SPARQL que se trouve la réponse à ce problème : celui-ci contient l'ensemble des variables requêtées, que des résultats y soient associés ou non (voir la clé `vars` dans le `head` d'une réponse SPARQL : 13). Pour construire la base de connaissance à partir des réponses SPARQL, il suffit donc de récupérer ces variables dans l'en-tête, et ensuite d'extraire les valeurs qui y sont associées dans les différentes « solutions de requêtes » contenues dans le document. Il est cependant nécessaire de supprimer les redondances dues au calcul d'un produit cartésien

entre les résultats. Pour ce faire, une valeur n'est liée à un résultat que si elle n'a pas déjà été extraite.

Une difficulté supplémentaire apparaît si une erreur a lieu au lancement d'une requête. Dans ce cas, la requête est relancée, mais cette fois-ci, le format de réponse demandé est le **XML**, au lieu du **JSON**. Cela permet d'éviter d'éventuels problèmes dans la construction du document **JSON** ; relancer la requête permet aussi de faire une seconde tentative lorsque la première requête a dépassé la durée maximale d'exécution autorisée. Le script de traitement des réponses décrit ci-dessus ne fonctionne cependant pas avec des données en **XML**. Plutôt que de réécrire l'intégralité de ce script pour du **XML** (ce qui serait également possible), le choix a plutôt été fait de traduire une réponse **XML** en un document **JSON** conforme avec la spécification **SPARQL** <sup>9</sup>. Ce **JSON** a une structure qui permet la traduction des résultats **SPARQL** pour constituer la base de connaissances. Le processus de transformation du **XML** en **JSON** en lui-même est assez simple, puisque les deux formats de réponse ont la même structure, mais simplement une syntaxe formelle différente. Ainsi, plutôt que de traduire directement la réponse **SPARQL XML** obtenue au format voulu pour la base de connaissance, il est plus facile de passer par un format intermédiaire en **JSON**.

Une fois ce traitement réalisé pour les quatre requêtes faites par entité, il ne reste plus qu'à ajouter à la base de connaissance ces données récupérées de *Wikidata*.

---

9. *Ibid.*

# Chapitre 7

## Un corpus augmenté : mise en relation du corpus aux données de *Wikidata* et possibilités ouvertes par cet alignement

### 7.1 Lier la TEI aux données nouvellement produites

Jusqu’alors, les catalogues ne servent que de sources de données qui permettent de réaliser l’alignement : ils n’ont pas été modifiés durant tout ce processus. Cependant, d’importantes modifications ont eu lieu pour l’ensemble du projet : les auteur.ice.s de manuscrits ont été liés à entités *Wikidata*, ce qui a permis la constitution d’une base de données « sur mesure ». Celle-ci contient des données pertinentes pour mener une étude économétrique qui cherche à faire le lien entre des facteurs biographiques et le prix d’un manuscrit ; à l’aide de ces données, il est possible de mieux évaluer la célébrité d’un.e auteur.ice, voire peut-être de comprendre qu’est-ce qui fait cette célébrité. Comme nous allons le voir plus bas, cette base de données peut être utilisée au delà de cette étude économétrique.

La première étape, cependant, est d’explicitier la relation entre les catalogues et cette source de données externe. Il est donc nécessaire d’annoter le corpus de catalogues à l’aide des identifiants *Wikidata*. Grâce son usage des attributs, la TEI rend possible l’annotation de corpus de façon assez intéressante. Un attribut permet de rajouter ou de normaliser des informations sur un élément sans en modifier le contenu textuel. Là où le corps du texte est contenu entre une balise ouvrante et une fermante, un attribut se trouve dans la balise ouvrante. Les attributs permettent, entres autres, des liens entre différentes parties du document, ou, ce qui nous intéresse dans ce contexte, avec des données extérieures. Il existe plusieurs manières, plus ou moins complexes, de créer des relations entre des

éléments d'un texte et des ressources externes avec la TEI. En suivant les préconisations des *Guidelines*<sup>1</sup>, les identifiants *Wikidata* ont été inclus dans un attribut `ref`. Cet attribut permet de pointer vers une ressource externe. Pour tirer véritablement parti d'un encodage numérique, l'alignement entre entités *Wikidata* et auteur.ice.s de manuscrits va un peu plus loin. Dans les catalogues mis à jour, ces identifiants sont précédés du préfixe « `wd:` » : `<name type="author" ref="wd:Q231457">Marie-Amélie</name>`. Celui-ci est une abréviation de *Wikidata*. Il permet à un.e lecteur.ice humain.e de comprendre que l'attribut `ref` renvoie à cette base de connaissance ; mais il permet surtout de créer des liens dynamiques vers celle-ci. En effet, ce préfixe est déclaré dans l'en-tête des catalogues (le `tei:teiHeader`) à l'intérieur d'un élément `tei:prefixDef` (16). Là, il est indiqué pour les lecteur.ice.s que le préfixe renvoie à *Wikidata*. Mais cet élément est surtout intéressant lors du traitement du document par une machine. Il y est indiqué que ce préfixe peut être remplacé par l'URL `https://www.wikidata.org/wiki/`. Ainsi, un processeur peut remplacer tous les préfixes par cet URL, et ainsi reconstruire des liens vers les pages des entités sur *Wikidata*. Cette manière de lier un `tei:name` à un identifiant *Wikidata* permet d'explicitement les liens entre les catalogues et la base de données propre au projet constituée à l'aide de SPARQL ; elle permet également de faire des renvois dynamiques vers des ressources en ligne, facilitant l'accès d'utilisateur.ice.s à celles-ci. Cette double méthode de renvoi permet à la fois un meilleur traitement par une machine et une lecture « augmentée » grâce à des liens hypertextes. La transformation des identifiants en URL complets peut également être explicitement mise en place dans des étapes de traitement ultérieures des catalogues. Il serait par exemple possible d'intégrer ces hyperliens au site Web de *Katabase*. Ainsi, le traitement automatisé des catalogues et l'alignement avec *Wikidata* va dans la continuité des objectifs de la TEI : permettre la création de documents enrichis pour des lecteur.ice.s et manipulables par des machines.

## 7.2 Que faire des entités ?

Le processus décrit dans cette partie d'alignement des `tei:name` avec des entités *Wikidata* et de construction d'une base de connaissances a été commencé avec une visée précise. Récupérer une grande quantité de données biographiques sur les auteur.ice permettrait de comprendre ce qui fait le prix d'un manuscrit. Le liage et la récupération de données via SPARQL ayant pris un temps important, cette étude n'a malheureusement pas pu être menée à son terme.

Au vu du temps passé à produire ces données, il serait cependant dommage de limiter l'usage des données obtenues à cette étude. En effet, le liage d'entités nommées ouvre énor-

---

1. TEI Consortium, *P5 : Guidelines for Electronic Text Encoding and Interchange*, Text Encoding Initiative, Version 4.4.0, 2022, URL : `https://tei-c.org/release/doc/tei-p5-doc/en/html/index.html`, 16. Linking, Segmentation and Alignment.

mément de possibilités : il permet d’obtenir des données structurées et manipulables par une machine, et donc de réaliser un traitement et une analyse « de masse » de documents textuels. Comme cela a déjà été dit, le liage d’entités a souvent, en humanités numériques, des visées cartographiques. Cette cartographie est souvent – et logiquement – consacrée à des lieux physiques<sup>2</sup>. *Bentham Project* a proposé une approche intéressante : une cartographie d’un corpus du philosophe Jeremy Bentham (philosophe utilitariste célèbre pour ses travaux sur le panoptique) encodé en XML-TEI. Les entités traitées par le projet Bentham dépassent les personnes, puisque des concepts sont extraits et liés à *Wikidata* ; elles sont identifiées par comparaison à des concepts présents sur la base en ligne *DBpedia* et en extrayant les thèmes principaux du corpus<sup>3</sup>. L’application Web conçue pour le corpus tire partie de ce liage d’entités : les concepts représentés dans des réseaux interactifs navigables ; les concepts majeurs dans la production de J. Bentham pour chaque décennie sont également visibles sous la forme de cartes de fréquentation (*heatmaps*)<sup>4</sup>. Ce projet – fruit d’une transcription campagne de transcription participative (*crowdsourcing*)<sup>5</sup> sur plus de huit ans – est d’une ampleur sans commune mesure à celle de *MSS / Katabase*. Ce projet offre tout de même des pistes intéressantes pour le développement Web dans un projet d’humanités numériques (ce qui constitue la troisième partie du présent mémoire). Visualiser les données extraites de *Wikidata* sur l’application Web de *Katabase* serait une manière intéressante de rendre les données extraites disponibles au public.

Au delà de questions cartographiques, la base de connaissance constituée ouvre la voie à d’autres analyses du corpus : elle contient des informations qui ne sont pas explicitement présentes dans le corpus de catalogues. L’extraction d’informations de *Wikidata* permettrait de procéder à la reconnaissance d’entités nommées dans les parties les moins structurées des catalogues (c’est-à-dire, les `tei:note` et les `tei:desc`, souvent riches d’informations sur les destinataires des lettres) : il serait possible, uniquement par détection de motifs, d’identifier les destinataires de lettres qui sont également auteur.ice.s de manuscrits présents dans le corpus. La constitution d’une base de connaissance permettrait alors de continuer à développer des méthodes de reconnaissance et de liage d’entités nommées à l’aide de techniques simples. Plutôt que d’utiliser l’apprentissage machine (pour lequel des méthodes sont déjà mises en place et pratiquées), coûteux en ressources, la présence de

---

2. A. Soudani, Y. Meherzi, A. Bouhafs, *et al.*, « Adaptation et évaluation de systèmes de reconnaissance et de résolution des entités nommées pour le cas de textes littéraires français du 19<sup>ème</sup> siècle »... ; Francesca Frontini, C. Brando, M. Riguet, *et al.*, « Annotation of Toponyms in TEI Digital Literary Editions and Linking to the Web of Data »... ; Noémie Boeglin, Michel Depeyre, Thierry Joliveau et Yves-François Le Lay, « Pour une cartographie romanesque de Paris au XIX<sup>e</sup> siècle. Proposition méthodologique », dans *Conférence Spatial Analysis and GEOMatics*, Nice, France, 2016 (Actes de la conférence SAGEO’2016 - Spatial Analysis and GEOMatics), URL : <https://hal.archives-ouvertes.fr/hal-01619600>.

3. Pablo Ruiz Fabo et Thierry Poibeau, « Mapping the Bentham Corpus : Concept-based Navigation », *Journal of Data Mining & Digital Humanities*, Atelier Digit\_Hum (Data deluge : which skills for...[ 2019]), p. 5044, DOI : 10.46298/jdmh.5044, p. 8-11.

4. *Ibid.*, p.14-17.

5. *Ibid.*, p. 2.

données structurées pourrait de continuer à développer des approches de « basse technologie » pour analyser un corpus semi-structuré. Grâce à cette base, il devient par également possible d'étudier la répartition des genres des auteur.ice.s de manuscrits : leur genre a été récupéré dans des requêtes SPARQL. Il a été envisagé de mener une telle étude au début de mon stage. Il est cependant difficile de faire cette analyse en s'appuyant uniquement sur les catalogues : le genre n'y est pas explicitement donné. Il est difficile d'appliquer des méthodes de détection de motifs à la question du genre d'un.e l'auteur.ice. Il serait possible de détecter les adjectifs au féminin dans le `tei:trait` ; cependant, de telles informations manquent souvent des entrées de catalogue. Une identification du genre à partir des prénoms serait presque irréalisable, du fait de la variété des prénoms, de la présence de prénoms non genrés dans la langue française, ou encore de l'absence de règles précises distinguant un prénom masculin et féminin. Le liage d'entités avec *Wikidata* permet, lui, de mener de telles études. Le genre d'un.e auteur.ice pourrait être annoté dans les corpus, comme cela a été fait pour les identifiants *Wikidata*. Étant donné que les identifiants explicitent le genre, ils permettent également la constitution d'une base de connaissances contenant des noms et prénoms masculins et féminins ; il serait alors possible d'identifier la présence de ces prénoms dans le texte afin d'annoter automatiquement des corpus de catalogue avec ces informations.

Enfin, l'alignement des auteur.ice.s de manuscrits peut permettre des enrichissements ultérieurs : plusieurs identifiants ont été récupérés en plus de l'identifiant *Wikidata* (« VIAF, DataBnF »...). Cela facilite les enrichissements ultérieurs de la base de connaissances constituée lors de cette étape en réalisant d'autres requêtes. Il serait également possible d'enrichir les contenus diffusés au public via l'application Web *Katabase*, en créant par exemple des pages pour les auteur.ice.s. Jusqu'à maintenant, cela n'a pas été fait car les manuscrits sont pauvres en informations biographiques. Ce cas de figure (une grande quantité d'entrées contenant peu d'informations) est similaire à celui présenté par Aguirre et al.<sup>6</sup> au sujet des collections de *Digitised cultural heritage* (Partimoine culturel numérisé). Ces collections, composées de millions de cartels d'œuvres encodées en XML en suivant le standard « Dublin Core », contient très peu d'informations sur chaque item. Leur proposition est de développer une chaîne de traitement sur le modèle du liage d'entités, afin de lier les collections à des pages *Wikipedia*. Grâce à l'alignement avec *Wikidata* déjà réalisé dans le cadre du projet *Katabase*, il n'est pas nécessaire de mener une nouvelle campagne de résolution d'entités. À partir d'une simple requête SPARQL, les URL des pages *Wikipedia* liées aux auteur.ice.s peuvent être automatiquement récupérées. Avec ces URL, il est facile d'automatiser la récupération du texte d'articles de *Wikipedia* pour

---

6. Eneko Agirre, Ander Barrena, Oier Lopez de Lacalle, Aitor Soroa, Samuel Fernando et Mark Stevenson, « Matching Cultural Heritage items to Wikipedia », dans *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, dir. Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk et Stelios Piperidis, Istanbul, Turkey, 2012, p. 1.

enrichir le site Web de *Katabase*.

## 7.3 En conclusion

Le processus d'alignement avec *Wikidata* présenté dans cette partie propose de s'appuyer sur des méthodes de détection de motif à des fins de NEL. Ce choix technique est simple, par opposition aux outils d'apprentissage machine généralement utilisées à cette fin. Grâce à la nature semi-structurée du corpus, il est cependant possible d'identifier très précisément les éléments pertinents pour l'interaction avec *Wikidata*. Il est également possible d'attribuer à ces éléments une valeur sémantique, ce qui permet de construire un algorithme de requêtes réactif. L'utilisation de technologies « simples » a cependant été contrebalancé par la conception d'algorithmes complexes, qui permettent de s'adapter de très près aux documents sources. La mise au point d'un algorithme construisant et lançant des requêtes sur le moteur de recherche de *Wikidata* contraste également avec les processus habituellement utilisés pour le liage d'entités nommées. Dans de tels contextes, la tendance est de comparer une chaîne de caractères dans le texte (reconnue comme une entité) avec une entité dans une base de connaissances et, à l'aide d'un algorithme, de procéder à un alignement. Cet algorithme peut, comme dans le cas de REDEN, s'appuyer sur des graphes qui représentent les candidats possibles pour un alignement ; l'entité de la base de connaissance choisie pour l'alignement est l'élément central de ce graphe<sup>7</sup>. Cependant, ces différents algorithmes développent leurs propres méthodes d'alignement et de sélection des candidats ; celles-ci ne s'appuient pas sur l'utilisation de moteurs de recherche déjà existants et performants ; elles ne semblent pas non plus prendre en compte les informations biographiques qui ont été utilisées ici. À l'inverse, l'alignement par détection de motifs s'appuie fortement sur l'utilisation d'un moteur de recherche, et met donc à profit les outils déjà mis au point par les bases de connaissance. De plus, la sélection du candidat présentée ici s'appuie sur l'extraction d'un grand nombre d'informations biographiques de la source, qui ne sont pas utilisées en NEL. Dans une démarche d'alignement d'entités nommées, s'appuyer au maximum sur les données disponibles dans le texte source peut être une perspective intéressante. En maximisant la quantité d'informations extraites, il est possible de réduire le bruit et de sélectionner le bon candidat avec une plus grande certitude. Le principal problème des méthodes d'alignement avec *Wikidata* présentées ici est qu'elles sont difficilement portables sur d'autres données : les motifs repérés, les tables utilisées pour convertir des données et les algorithmes d'extraction ne fonctionnent qu'avec les catalogues de ventes de manuscrits. L'utilisation de l'API *Wikidata* est, elle, plus facilement adaptable. À l'inverse, l'apprentissage machine offre la promesse d'une plus grande portabilité ; il est cependant nécessaire d'entraîner des modèles afin de les adapter aux données sélectionnées.

---

7. C. Brando, Francesca Frontini et J.G. Ganascia, « REDEN... », p. 72.

Il est intéressant de remarquer que les questions soulevées lors de cette étape se rapprochent des problématiques propres à l'apprentissage machine. Le premier problème apparent est un possible biais dans le processus d'extraction d'informations. Celui-ci est le fruit d'une analyse du corpus, d'une manipulation du moteur de recherche de *Wikidata* et de nombreux tests ; cependant, il est possible que certains aspects du corpus aient été par erreur favorisés. Ce problème se retrouve dans les données de test utilisées : puisque ces données ont fait l'objet de plus d'attention que le reste du corpus, il est possible que le processus d'extraction d'informations soit plus performant sur celles-ci que pour le reste du corpus. Cette difficulté rappelle le « surapprentissage » (*overfitting*) fréquent en apprentissage machine : un algorithme « apprend » trop efficacement des données sur lesquelles il est entraîné ; ses performances deviennent alors médiocres sur d'autres données, puisqu'il ne sait pas s'y adapter. Pour mesurer efficacement ce risque, deux jeux de données de test sont utilisés : le jeu d'entraînement (que l'algorithme utilise pour « apprendre ») et le jeu d'évaluation (qui mesure les performances). Ici, un seul jeu de données est utilisé, puisqu'il n'y a pas à proprement parler de phase d'entraînement. Il serait cependant utile de créer un jeu d'évaluation distinct afin d'évaluer les performances du liage d'entités. Cela permettrait de donner des résultats plus adéquats en évitant que l'algorithme soit évalué sur des données auxquelles il est trop adapté. En plus de ces problèmes liés à la source de données (ici, les catalogues), d'autres problèmes viennent de la base de connaissance utilisée dans la résolution d'entités. Comme le remarquent Carmen Brando et al.<sup>8</sup>, la capacité à aligner une entité dans un texte avec la bonne entité sur une base de connaissance dépend largement des données présentes dans celle-ci. L'extraction d'informations et l'utilisation du moteur de recherche de *Wikidata* ont cherché à contourner ce problème en traduisant les données sources en leurs équivalent sur *Wikidata*. Cela permet de se conformer aux biais de la base de connaissance, ce qui n'est pas possible avec des méthodes traditionnelles d'apprentissage machine : lorsqu'une entité est résolue à partir de sa centralité dans un graphe, la différence de vocabulaires dans les sources et la base de connaissance ne peut pas être prise en compte.

Il va de soit que la détection de motifs n'est pas toujours possible : si elle fait sens pour un corpus semi-structuré, elle serait probablement moins efficace avec un texte littéraire, où les structures des phrases sont moins codifiées. L'encodage en TEI des sources textuelles peut, cependant, faciliter l'utilisation de telles méthodes sur des corpus plus complexes. En effet, puisqu'un tel encodage permet un balisage sémantique, il est possible de cibler des éléments précis à l'intérieur desquels détecter des motifs. Ainsi, il n'est pas nécessaire de trouver des méthodes qui fonctionnent pour l'intégralité du texte, et de travailler à une échelle restreinte.

Cette phase d'enrichissement à l'aide de données externes ouvre la voie à de nouvelles

---

8. C. Brando, N. Abadie et Fontini Frontini, « Évaluation de la qualité des sources du web de données pour la résolution d'entités nommées »..., p. 37-8.



études, et donc à une meilleure compréhension du corpus ; elle pourrait aussi permettre d'enrichir l'application *Katabase* avec de nouvelles fonctionnalités (des pages personnelles sur les auteur.ice.s de manuscrits, par exemple). Cependant, traiter et analyser des données n'est pas la seule manière de comprendre le corpus. Pour mieux appréhender celui-ci, y compris dans une approche économique, la visualisation des données a également un intérêt. Elle permet une représentation synthétique et ciblée d'informations, et offre donc une interprétation du corpus, qui, comme nous le verrons dans la partie suivante, permet de mieux le connaître. Les questions de visualisation ne sont cependant pas que des questions d'ingénierie de la donnée : ce sont également des questions graphiques et esthétiques. Dans un projet d'humanités numériques, les questions graphiques gagnent en importance lorsque les résultats de la recherche doivent être diffusés à un public par l'intermédiaire d'une application Web.



## Troisième partie

Rendre la recherche réutilisable et  
interopérable : *KatAPI*, une API  
pour échanger des données  
structurées



La problématique de cette partie est la suivante : comment rendre la recherche en humanités numériques réutilisable et encourager le partage de données entre projets de recherche ? La réponse à cette question, au sein du projet *MSS / Katabase* prend la forme d'une API. Celle-ci diffuse des données source du projet, mais aussi des informations issues de la recherche et du traitement computationnel de celles-ci dans deux formats structurés : le JSON et le XML-TEI.

En guise d'introduction, il me semble important ici de faire un bref rappel des principes fondamentaux du développement applicatif et de l'architecture du Web. Certains de ces principes ont déjà été présentés, mais ce rappel permettra de rendre la suite de ce chapitre plus claire.

Le Web est une architecture de type client-serveur qui permet l'interaction de machines dans un réseau décentralisé. Cette architecture définit deux fonctions possibles pour une machine : client (qui envoie et récupère des données distantes) et serveur (qui contient une application permettant l'interaction avec le client, ainsi qu'une base de données). Dans le Web, l'architecture client-serveur repose sur un protocole technique qui pose les principes de l'interaction entre des machines : le HTTP. Selon ce principe, l'accès d'un client à une ressource stockée sur un serveur (une page Web, par exemple) se fait au travers d'un URL. Celui-ci sert à localiser la ressource. L'interaction est définie par des requêtes (du client au serveur) et par des réponses (du serveur au client). Six types de requêtes sont possibles, qui correspondent à autant d'actions : demander des données à un serveur, en supprimer, les mettre à jour... Requêtes et réponses ont la même structure : elles sont composées d'un en-tête HTTP (qui contient des métadonnées sur la requête ou la réponse) et d'un corps (qui contient la requête ou la réponse en elle-même ; celui-ci peut être vide). Parmi les métadonnées contenues dans l'en-tête, deux sont importantes dans notre contexte (surtout pour les réponses du client au serveur) : le statut et le type MIME. Un statut est un code à trois chiffres qui caractérise la réponse et le déroulement de l'interaction client-serveur. « 200 » indique que tout va bien, « 404 » que la ressource demandée n'existe pas (c'est-à-dire, que l'URL de la requête ne pointe pas vers une ressource existante). Le type MIME est une manière standardisée de définir le format d'une réponse : « `application/json` » lorsque la réponse est au format JSON, « `application/xml` » lorsque du XML est renvoyé.

Une API est une application, c'est-à-dire un ensemble de fonctions rassemblées pour permettre l'interaction entre un client et un serveur. Souvent, et c'est le cas pour les applications du projet *MSS / Katabase*, une application permet l'interaction avec une base de données, elle-aussi stockée sur le serveur. Par exemple, lorsque l'on accède à une page *Wikipedia*, l'application reçoit une requête de la part d'un client ; elle va chercher dans la base les données pertinentes, les structure sous la forme d'une page HTML et les renvoie à l'utilisateur, qui verra la page s'afficher sur son écran. La principale différence entre une API Web et une application (comme *Wikipedia*) est que l'API est « un site Web pour

machines ». Elle ne contient pas d'interface graphique et les réponses renvoyées par l'API sont des données dites « brutes » ; celles-ci ont l'avantage d'être lisibles et manipulables par une machine. L'interaction avec une API se fait en envoyant une requête à l'application avec un URL qui contient les paramètres de la requête et des valeurs associées à ces paramètres (quelles données demander, quel format de réponse utiliser...) et éventuellement un corps de requête qui contienne des données (ce qui n'est pas le cas lorsque la requête correspond à une demande de données). L'URL contient donc la sémantique de la requête, il correspond à une question posée à un serveur de façon structurée. L'intérêt d'une telle API est que, pour peu que l'on sache interagir avec elle, il est possible de récupérer automatiquement des données structurées, prêtes à être utilisées. Pour télécharger un jeu de données, il n'y a pas besoin, par exemple, d'aller sur la bonne page Web, de cliquer sur un bouton « Télécharger » et d'enregistrer le fichier au bon endroit. Toutes ces opérations peuvent être faites automatiquement, comme cela a été fait pour plus de 80000 requêtes lors de l'interaction avec *Wikidata* décrite précédemment. Cette perspective est particulièrement intéressante pour un projet de recherche, puisqu'elle peut permettre de récupérer automatiquement les jeux de données produits par un autre projet, ce qui encourage la réutilisation des données de la recherche. Pour des ingénieur.e.s des données, les API facilitent grandement le travail de récupération automatisée d'informations : elles permettent de « se concentrer sur le quoi, et non sur le comment »<sup>9</sup>.

Souvent, un site Web peut être composé d'une version accessible pour des utilisateur.e.s humain.e.s et pour des machines ; c'est le cas dans le projet *Katadbse*. L'application pour humain.e.s a déjà été développée ; c'est la version pour les machines, son fonctionnement et son utilité pour permettre la réutilisation des données du projet qui est donc présentée dans cette partie.

---

9. *focus on the what, and not on the how*. Citation issue de Jaimie Murdock et Colin Allen, « InPhO for All : Why APIs Matter », *Journal of the Chicago Colloquium on Digital Humanities and Computer Scienc*, 1–3 (2011), URL : <http://www.jamram.net/docs/jdhcs11-paper.pdf> (visité le 12/07/2022), p. 2.

## Chapitre 8

# Standards de design et statut des API dans pour les humanités numériques centrées sur le texte

Il n'existe aucune norme communément acceptée pour la conception d'une API : dans l'absolu, tant que celle-ci fonctionne, elle peut être considérée comme une API. Cependant, des standards se sont développés afin de regrouper les différentes méthodes et de développer des bonnes pratiques de design. Deux types de standards sont ici présentés : les standards d'architecture et les standards d'interopérabilité. Un standard d'architecture ou de design, comme le REST, définit comment concevoir l'interaction client-serveur dans une API et proposent une manière de structurer les réponses du serveur. Il ne définit cependant pas quels paramètres de recherches sont autorisés, ni le corps d'une réponse doit être organisé. Il ne définit pas non plus précisément comment organiser une requête, c'est-à-dire quel URL construire pour accéder aux données. C'est le rôle des standards d'interopérabilité (*Open Archives Initiative Protocol for Metadata Harvesting* (OAI-PMH) , *Canonical Text Services* (CTS) , *Distributed Text Services* (DTS) ). Ceux-ci définissent la sémantique des requêtes (c'est-à-dire les paramètres autorisés) et la manière de construire un URL pour accéder à certaines données. Ils permettent donc l'interopérabilité, puisqu'une même syntaxe peut être utilisée pour requêter différentes API. Ni le REST, ni les standards d'interopérabilité ne définissent l'implémentation d'une API : les deux types de standards s'intéressent à l'interaction du client avec l'API, ils ne définissent pas l'application en elle-même. Une application qui se conforme à un standard garantit simplement un certain type d'interaction client-serveur. En plus de ces deux standards qui sont ici présentés, les principes *Findable Accessible Interoperable Reusable* (FAIR) sont également décrits : ceux-ci ne définissent ni implémentation, ni sémantique, mais sont simplement des principes généraux pour la manière dont des données doivent être partagées dans un contexte de science ouverte et réutilisable.

## 8.1 Que faire du FAIR ? Le partage des données de la recherche

Les principes FAIR<sup>1</sup> sont essentiels aux problématiques de la science ouverte. Celle-ci encourage non seulement la diffusion de la recherche dans des formats traditionnels (articles, conférences...), mais aussi la diffusion de données brutes issues de la recherche scientifique. L'objectif, dans la diffusion de ces données brutes, est de rendre les données réutilisables, mais aussi de permettre la reproductibilité de la recherche, et donc la vérification des résultats obtenus.

La problématique de la diffusion de données recherche brutes se retrouve bien sûr en humanités numériques. Dans le domaine de la reconnaissance optique de caractères, par exemple, se développe le partage de vérités de terrain. La reconnaissance de caractères est permise par l'entraînement d'un algorithme d'apprentissage machine sur des textes qui ont été transcrits manuellement ; c'est ce couple texte numérisé/transcription qui constitue la vérité de terrain. Or la production de ces données demande beaucoup de temps ainsi que des compétences spécifiques pour des graphies difficiles à lire et des langues rares<sup>2</sup>. Pour faciliter le développement de méthodes de reconnaissance optique de caractères se développent donc des initiatives pour faciliter le partage de vérités de terrains. Parmi ces initiatives se trouve le projet *HTR United*<sup>3</sup>, qui offre un portail pour l'accès à et le dépôt de vérités de terrains. Celles-ci sont partagées avec des métadonnées précisant par exemple leur contexte de production, le domaine auquel appartient au corpus.

La diffusion de données de recherche est cependant une opération plus complexe qu'il ne pourrait y paraître ; c'est pourquoi les principes FAIR ont été édités en 2014<sup>4</sup>. Ils forment une base minimale de recommandations sur la manière de partager des données de recherche. Cette base minimale est composée de quatre principes<sup>5</sup> :

- *Findability* : les jeux de données doivent être diffusés sur des dépôts publics ; ils doivent être décrits avec suffisamment de métadonnées pour que leur périmètre et la manière dont ils ont été constitués soient clairement identifiables.

---

1. Martin Boeckhout, Gerhard A. Zielhuis et Annelien L. Bredenoord, « The FAIR guiding principles for data stewardship : fair enough ? », *European Journal of Human Genetics*, 26–7 (juil. 2018), p. 931-936, DOI : 10.1038/s41431-018-0160-0.

2. Un exemple de ces corpus « difficiles » est le développement de la reconnaissance optique de caractères pour le Cham, langue parlée par une minorité originaire du Sud-Est de l'Asie. Celle-ci comprend environ 30000 locuteurs. Le projet travaille également à la reconnaissance du Cham ancien, qui est extrêmement méconnu (Anne-Valérie Schweyer, Jean-Christophe Burie et Tien Nam Nguyen, « Analyse, Reconnaissance et Indexation des manuscrits CHAM », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris, 2022 ).

3. Alix Chagué et Thibault Clérice, « Sharing HTR datasets with standardized metadata : the HTR-United initiative », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris, France, 2022, URL : <https://hal.inria.fr/hal-03703989>.

4. M. Boeckhout, G. A. Zielhuis et A. L. Bredenoord, « The FAIR guiding principles for data stewardship... », p. 931.

5. *Ibid.*, p. 932-933.



- *Accessible* : les jeux de données doivent être accessibles librement et idéalement de façon manuelle aussi bien qu'automatique.
- *Reusable* : les jeux de données doivent être réutilisables. Cela implique de définir un périmètre pour la réutilisation des données, notamment au travers de licences.
- *Interoperable* : les jeux de données doivent être diffusés dans des formats structurés, compréhensible par des humain.e.s et, idéalement, manipulables par des machines.

Qu'en est-il de la relation entre ces principes et la conception d'API ? La diffusion automatisée par le biais d'applications semble aller main dans la main avec les principes FAIR. Une API permet de diffuser automatiquement des données brutes et structurées dans des formats interopérables. Cependant, une API n'adhère pas par défaut aux principes FAIR. Elle garantit l'accessibilité des données ; pour qu'elles soient réutilisables, la conception d'API doit être considérablement étoffée. Il faut en effet définir le contexte de production des données, ainsi que la licence dans laquelle elles sont diffusées. Il est donc nécessaire d'ajouter des métadonnées aux réponses de l'API ou de trouver un moyen de documenter les conditions de réutilisation des données. La notion de « dépôt public » requis par les principes FAIR pour partager des données change également de sens. L'API sert pas à partager les données sur un autre dépôt qui agrège des jeux de données, elle sert à les diffuser à des utilisateur.ice.s. Une API pour la recherche forme donc sa propre plateforme de diffusion des données. Enfin, il faut garantir que les données soient compréhensibles. Cet aspect des principes FAIR n'est pas anodine dans la conception d'une API. Non seulement celle-ci doit renvoyer des réponses structurées et valides ; la sémantique des réponses doit également être explicitée. Cela veut dire que, soit dans l'API, soit dans la documentation, le sens de chaque élément doit être précisé. La documentation est ici importante ; mais l'utilisation de la TEI est également intéressante en vue de ce besoin de produire des données compréhensibles. Ce standard, en plus d'être très pratiqué dans les humanités numériques, est défini dans les *TEI Guidelines*<sup>6</sup> et dispose donc d'une documentation officielle. Au delà de ces principes abstraits, de véritables standards techniques peuvent servir d'inspiration dans le développement d'API pour diffuser des données de recherche.

## 8.2 Le REST : un modèle pour le design d'API

Ici est présenté le standard REST, défini par Roy Fielding<sup>7</sup> ; ce standard architectural est considéré comme l'idéal à atteindre en matière de design d'API. Il n'a cependant pas été suivi ici, en partie parce qu'il est mieux adapté à un projet de plus grande échelle, et en partie car il ne fonctionne pas de façon optimale avec le modèle de données défini,

---

6. TEI Consortium, *P5 : Guidelines for Electronic Text Encoding and Interchange...*

7. R. Fielding, *Architectural Styles and the Design of Network-based Software Architectures...*

et notamment avec le choix de construire une API renvoyant des réponses conformes à la TEI.

Le standard REST (« État de transfert représentatif ») est un standard et un style pour la conception d'API dans une architecture client-serveur. Il est aujourd'hui un standard de fait pour la conception de tels outils, et est considéré comme une bonne pratique dans l'architecture d'API. Le standard REST ne définit pas comment une API doit être implémentée, mais définit la structure de l'interaction client-serveur. Historiquement, le développement de ce standard est fortement lié à l'établissement du protocole HTTP, puisque le concepteur du REST est l'un des ingénieur.e.s à l'origine du protocole. Le REST définit un ensemble de quatre contraintes supplémentaires à celles présentes dans le protocole HTTP<sup>8</sup>. Une API REST doit répondre aux contraintes suivantes :

- Être sans état<sup>9</sup>. Cela signifie qu'une requête du client au serveur doit contenir l'ensemble des informations nécessaires au traitement de la requête. Le serveur ne peut accéder à un contexte ou à d'autres informations pour traiter la requête. Cela permet de rendre les requêtes plus faciles à interpréter, limite le risque d'erreurs et permet à l'API de fonctionner à plus grande échelle, puisqu'elle n'a pas besoin d'utiliser de les ressources pour chercher des informations ailleurs que dans la requête.
- Présenter une interface uniforme<sup>10</sup>. C'est peut-être ce principe qui est définitoire du REST : l'interface, c'est-à-dire la partie de l'API qui est exposée au client, doit toujours être la même. Il y a donc un découplage entre la communication avec le client (réception des requêtes et envoi des réponse) et le traitement des requêtes, qui ne sont pas accessibles au client. Cette séparation permet de modifier la manière dont une requête est traitée sans devoir changer la sémantique et la structure de l'interaction avec le client. Ce principe implique d'autres contraintes :
  - L'identification des ressources dans des requêtes : un client requête une URI et le serveur lui renvoie une représentation de cette URI. Dans le contexte du Web, cela veut dire qu'un URL doit nécessairement être fourni au moment de la requête et que celui-ci doit pointer vers les ressources à traiter.
  - Les ressources sont communiquées sous la forme de représentations. Cela veut dire que les données envoyées au client peuvent être dans une structure ou un format différents de ceux de la base de données : l'API doit traduire les données telles qu'elles sont dans la base en leur représentation.
  - Les réponses HTTP doivent être auto-descriptives : les en-têtes et le corps de la requête doivent contenir assez de (méta-)données pour pouvoir être manipulées par le client sans avoir besoin d'informations extérieures.

---

8. *Ibid.*, p. 94.

9. *Ibid.*, p. 78-79.

10. *Ibid.*, p. 81-82, 86-97.

- Utilisation de l'hypermédia comme moteur de navigation. Comme sur un site Web normal, des hyperliens doivent permettre de naviguer d'une ressource à l'autre et doivent être accessibles directement depuis l'API.
- Autoriser la mise en cache explicite<sup>11</sup>. Une réponse du serveur au client doit explicitement indiquer si les données contenues dans une réponse peuvent être mises en cache (c'est-à-dire, stockées dans la mémoire pour y être accédées plus tard). La possibilité de cacher les données peut potentiellement diminuer le nombre de requêtes envoyées au serveur.
- Présenter une architecture en couches<sup>12</sup>. Une application peut être entièrement stockée sur un serveur, ou fonctionner par l'interaction de différents serveurs. Dans ce cas, un client peut être connectés à différents serveurs, mais ne doit pas savoir s'il se trouve sur un serveur intermédiaire ou un serveur final. Les serveurs intermédiaires peuvent modifier les messages à renvoyer au client, répartir les charges sur différents serveurs ou établir des politiques de sécurité.
- Diffuser du code à la demande (optionnel)<sup>13</sup>.

Les concepts centraux dans ce système architectural – du moins, pour une petite API comme celle de *Katabase* sont le fait d'être sans état et l'uniformité de l'interface. Ces deux concepts signifient que les requêtes comme les réponses doivent être entièrement auto-descriptives. Les autres principes sont mieux adaptés à des API à plus grande échelle : la mise en cache explicite suppose que la même requête puisse être envoyée en série par un client à un serveur ; pour une API d'un projet de recherche, il est plus probable qu'une requête soit faite une unique fois pour récupérer des données. L'utilisation d'une architecture en couche n'est pas du ressort de *KatAPI*, puisque l'intégralité de l'API est présente sur un seul serveur. Enfin, la diffusion de code exécutable n'est pas non plus adaptée aux données qui sont requêtées et diffusées via l'API.

La mise au point d'une architecture sans état demande de mettre au point une sémantique pour les requêtes qui soit auto-descriptive, c'est-à-dire qui décrive l'intégralité de la requête. L'uniformité de l'interface, d'un point de vue de développement applicatif, est une problématique plus complexe. D'abord, elle demande d'établir des modèles de données uniformes pour le corps des réponses HTTP. Cela veut dire que, pour chaque format de réponse autorisé et point d'accès aux données possible, les réponses doivent être contenues dans une structure qui est toujours la même. Ensuite, les en-têtes doivent être constitués pour décrire adéquatement la réponse. Cela passe par la définition d'un statut HTTP et d'un type MIME à chaque réponse. Ainsi, une machine peut savoir quel format de réponse est obtenu, et comment la requête s'est déroulée. Cependant, l'utilisation d'hyperliens comme moteur de l'application n'a pas été implémentée : les réponses ne

---

11. *Ibid.*, p. 79-81.

12. *Ibid.*, p. 82-84.

13. *Ibid.*, p. 84-85.

contiennent pas d’hyperliens qui permettent de naviguer automatiquement vers d’autres ressources ou actions. L’ensemble des ressources accessibles et la manière d’y accéder sont au contraire définies de façon centralisée dans la documentation de l’API, qui n’est pas lisible par une machine. Enfin, du point de vue de l’implémentation, la communication client-serveur est indépendante du traitement des requêtes et de l’interaction avec les bases de données.

Par conséquent, les choix techniques mis en place font de l’API une application *RESTlike* (c’est-à-dire, qui se conforme à un certain degré aux principes REST), voire *RESTful* (qui se conforme entièrement aux principes REST). *KatAPI* suit le standard REST là où il est pertinent pour les données transmises (ce qui exclut la condition de transmission de code et la clause de mise en cache). L’API suit également les principes REST là où il est de son ressort de les suivre (ce qui exclut une architecture en couches, puisque l’application est distincte de son installation sur un serveur). Suivre ces principes et s’en inspirer permet de concevoir une application au comportement uniforme, et où le client n’a pas à vérifier les réponses obtenues, puisqu’elles seront toujours les mêmes et qu’elles contiendront toutes les informations nécessaires à leur traitement (format de réponse, statut HTTP...).

### 8.3 OAI-PMH, CTS et DTS : quels standards pour le partage du texte en humanités numériques ?

Les principes REST, présentés ci-dessus, concernent l’architecture et le design d’une API; cependant, ils ne s’attaquent pas à une autre partie du problème : l’interopérabilité des API. Chacune de ces applications définit sa propre sémantique pour les requêtes et autorise certains formats de réponse. Cela conduit à un problème pour les programmes consommateurs d’API : l’interaction avec deux API demande la construction de requêtes différentes, et demande donc d’apprendre à utiliser différentes API. Lorsque celles-ci sont complexes, apprendre à utiliser une API n’a rien d’anodin, et se servir de différentes API devient complexe. De plus, l’absence d’interopérabilité signifie qu’il est difficile, voire impossible, de travailler de la même manière avec les données issues d’API différentes. Si l’une renvoie du JSON mais l’autre du XML, par exemple, les résultats issus des deux API doivent être utilisés différemment. Cela complexifie la convergence des données de la recherche et limite l’interopérabilité de ces données, qui est pourtant l’un des principes FAIR<sup>14</sup>. Pour faire face à ce problème, des standards d’interopérabilité ont été développés. Ceux-ci définissent à minima une sémantique unique qui peut être implémentée par plusieurs API. Parfois, ils définissent également des formats de réponse. Le plus célèbre

---

14. M. Boeckhout, G. A. Zielhuis et A. L. Bredenoord, « The FAIR guiding principles for data stewardship... », p. 932-933.

de ces standards, dans les humanités, est le protocole IIIF, qui permet d'interagir avec une image (en lui ajoutant une couche d'annotations, par exemple). Dans les humanités centrées sur le texte, les standards d'API sont plus diversifiés. Trois sont présentés ici : l'OAI-PMH, le CTS et le DTS. Aucun d'entre eux n'a été adopté dans la conception de l'API *Katbase*, mais ils sont tous riches de leçons dans le développement d'une telle application

### 8.3.1 OAI-PMH : un premier standard à succès

L'OAI-PMH vient du monde des archives et des bibliothèques. Les bibliothèques ne conservant pas des objets uniques, mais des documents produits en série, elles sont pionnières dans le développement de méthodes computationnelles pour le partage des données<sup>15</sup>. Aussi ce standard a-t-il connu une adoption inégalée, puisqu'il est aujourd'hui le standard d'interopérabilité le plus utilisé dans les bibliothèques<sup>16</sup>; il est y compris utilisé par les plateformes en ligne HAL, OpenEditions ou encore CAIRN. Ce standard vise à permettre l'interopérabilité sémantique (uniformiser la sémantique des requêtes et des réponses) et technique (uniformiser les formats de réponse)<sup>17</sup>. Il ne vise cependant pas à partager des documents, mais seulement leurs métadonnées (en d'autres termes, l'OAI-PMH ne permet pas de partager un document, mais seulement sa notice bibliographique); il ne permet pas non plus de partager des données sur un document, mais sur une collection entière. Les réponses de l'OAI-PMH sont constituées de collections de notices encodées en XML; l'interopérabilité sémantique des réponses est permise par l'utilisation du Dublin Core, un format considéré comme un plancher pour le partage de données bibliographiques. L'interopérabilité technique, elle, repose sur un format de réponse en XML<sup>18</sup>. Comme le montrent Prime-Claverie et Mahé, cette interopérabilité n'empêche cependant pas une grande variété dans les implémentations : certaines utilisent des réponses en XML-TEI, d'autres non; les structures des réponses varient elles-aussi, avec différents éléments Dublin Core associés à différentes données<sup>19</sup>. D'un point de vue technique, le fonctionnement d'OAI-PMH est intéressant : l'entrepôt OAI-PMH est distinct de la base d'où il récupère les données. Cela veut dire que, lorsqu'une requête est reçue, l'API la communique à l'entrepôt OAI-PMH qui transforme les informations contenues dans la base de données. D'un point de vue d'ingénierie, cela veut dire que les parties de ce système (base de données – entrepôt OAI-PMH – API) sont distincts<sup>20</sup>. Par conséquent, les

---

15. Camille Prime-Claverie et Annaïg Mahé, « Le défi de l'interopérabilité entre plates-formes pour la construction de savoirs augmentés en sciences humaines et sociales », dans *Ecrilecture augmentée dans les communautés scientifiques*, 2017, URL : [https://archivesic.ccsd.cnrs.fr/sic\\_01511618](https://archivesic.ccsd.cnrs.fr/sic_01511618) (visité le 29/08/2022), p. 2-3.

16. *Ibid.*, p. 5.

17. *Ibid.*, p. 6-14.

18. *Ibid.*, p. 4.

19. *Ibid.*, p. 6-14.

20. *Ibid.*, p. 4.

données n'ont pas besoin d'être stockées dans des formats compatibles ; des transformations dans une partie n'impliquent pas de devoir transformer les autres parties (sauf pour permettre l'interaction entre les différents éléments du système). Par exemple, une modification de la structure de la base de données ne demande pas de transformer toute l'API, mais seulement la manière dont ces données sont transformées pour être compatibles avec le standard OAI-PMH.

Que retenir de cet exemple ? Tout d'abord, que ce n'est pas parce qu'un standard existe qu'il sera implémenté uniformément. Une marge de manœuvre est toujours présente, ce qui peut amener à des résultats très différents. Bien sûr, cela va à l'encontre du principe d'interopérabilité et complexifie la convergence des ressources d'origine différente. Cependant, comme le rappellent Prime-Claverie et Mahé<sup>21</sup>, ces variations ne sont pas nécessairement problématiques. Elles permettent aux différentes institutions qui l'implémentent de fournir des données adaptées à certaines communautés de pratique – il va de soit qu'une archive et une bibliothèque sont deux institutions très différentes, et qu'elles ne peuvent représenter leurs documents de la même manière. L'implémentation technique de l'OAI-PMH est également intéressante : elle sépare les données (dans une base de données) de leur représentation. Ce qui est communiqué au client, c'est cette représentation, compatible avec le protocole OAI-PMH. Cette idée, qui rappelle les principes REST, est mise en place d'autres standards présentés ici. D'un point de vue de développement, la séparation entre les différentes parties nécessaires au fonctionnement de cet API facilitent également son maintien et l'ajout de nouvelles fonctionnalités, puisque le système ne doit pas être modifié lorsque c'est seulement un élément qui est affecté. Cependant, l'OAI-PMH n'est pas adapté aux besoins de *Katabase* : ce standard ne permet pas d'échanger du texte, mais seulement des métadonnées. De par sa taille, il est également difficile à implémenter pour un projet de recherche.

### 8.3.2 CTS, DTS : des méthodes de standardisation pour l'échange de texte

Le CTS a été développé dans un contexte très différent de l'OAI-PMH. Il est issu d'un projet de recherche à Harvard, le *Homer Multitext* consacré à l'édition de l'*Illiad*e et de l'*Odyssée* d'Homère. Le standard est donc pensé pour l'édition de textes canoniques antiques. Il naît d'un problème technique : comment faire coïncider la structure en arbre d'une édition numérique avec la structure citationnelle traditionnelle, où les citations font références à des pages. Aux origines du projet se trouve également l'idée qu'un texte existe à l'intérieur d'un système citationnel, où les textes fonctionnent par référence les uns aux autres<sup>22</sup>. Cette idée est particulièrement importante dans un projet d'édition

21. *Ibid.*, p. 15-16.

22. David Neel Smith et Christopher W. Blackwell, « Four URLs, limitless apps : Separation of concerns in the Homer Multitext architecture », dans *Donum natalicium digitaliter confectum Gre- gorio*

critique, comme le *Homer Multitext*, où les différents témoins des textes d’Homère doivent pouvoir être alignés et comparés. La TEI a développé ses propres systèmes de référence pour les éditions critiques, basés sur le concept d’arbres décisionnels (où chaque témoin peut être représenté par une feuille, et où les différentes feuilles viennent d’une branche commune)<sup>23</sup>. Mais ce système ne permet pas de cibler efficacement un passage dans un texte, ni de faire référence à ce passage en dehors du document lui-même. Comment, dans ce cas, représenter des références dans une édition numérique ? Comment rendre des parties de texte accessibles à l’extérieur d’un texte, dans une architecture client-serveur ? La réponse proposée à l’initiative de Smith et Blackwell est le CTS. Ce système définit une sémantique permettant d’identifier un texte, un passage de texte ou même un ensemble de textes grâce à une URI. Celle-ci cible de plus en plus précisément le passage voulu. Par exemple, `urn:cts:greekLit:tlg0012.tgl001.msA` permet de cibler, à l’aide d’un identifiant CTS, le manuscrit `msA` du texte `tlg0012.tgl001` du corpus `greekLit`<sup>24</sup>. L’intérêt de ce système de pointeurs vers un texte est qu’il peut être adapté au Web, en précédant l’identifiant ci-dessus de l’URL d’une API supportant le standard CTS. Il est donc possible de cibler précisément et de récupérer un texte depuis un dépôt particulier. L’intérêt de ce système est également qu’il est généralisable à d’autres corpus ; le standard a donc été adapté par plusieurs API, qui ont dû définir leurs propres implémentations de ce standard : là où le standard est implémenté, il est possible d’accéder à l’aide d’une sémantique unique à n’importe quel(s) texte(s). Parmi ces implémentations peut être cité *CapiTainS*, développée pour la *Perseus Digital Library* et l’*Open Philology Project*<sup>25</sup>. CTS est donc intéressant d’un point de vue de l’ingénierie, puisqu’il permet la convergence des méthodes de partage et de diffusion du texte. Le standard est également intéressant d’un point de vue de la représentation du texte, puisqu’il permet de représenter un texte comme étant un graphe de citations et de références (une citation correspondant à une URI reliée à une autre par un prédicat). Ce système de représentation permet donc de s’éloigner de la structure hiérarchique des éditions numériques en TEI, une représentation critiquée depuis ses débuts parce qu’elle impose un ordre hiérarchique unique à un texte, pourtant irréductible à une seule interprétation<sup>26</sup>.

---

*Nagy septuagenario a discipulis collegis familiaribus oblatum : A Virtual Birthday Gift Presented to Gregory Nagy on Turning Seventy by His Students, Colleagues, and Friends*, 2012, URL : <https://chs.harvard.edu/d-n-smith-c-w-blackwell-four-urls-limitless-apps-separation-of-concerns-in-the-homer-multitext-architecture/> (visité le 22/08/2022).

23. TEI Consortium, *P5 : Guidelines for Electronic Text Encoding and Interchange...*

24. D. N. Smith et C. W. Blackwell, « Four URLs, limitless apps : Separation of concerns in the Homer Multitext architecture »...

25. Bridget Almas et T. Clérice, « Continuous Integration and Unit Testing of Digital Editions », *Digital Humanities Quarterly*, 11–4 (févr. 2018), Publisher : Alliance of Digital Humanities, URL : <https://hal.archives-ouvertes.fr/hal-01709868> (visité le 12/07/2022).

26. Allen Renear, Elli Mylonas et David Durand, « Refining our Notion of What Text Really Is : The Problem of Overlapping Hierarchies », dans *Research in Humanities Computing*, dir. Nancy Ide et Susan Hockey, Oxford, 1996, URL : <https://cds.library.brown.edu/resources/stg/monographs/ohco.html> (visité le 22/08/2022).

Le standard CTS n'est cependant pas parfait, et plusieurs critiques peuvent lui être faites. D'un point de vue de son développement, c'est un standard qui a été développé dans une certaine communauté de pratique, et qui n'est donc pas adaptable à tous les textes ; ensuite, c'est un format qui a été développé pour répondre aux besoins de chercheur.euse.s, et qui ne suit donc pas les meilleures pratiques en termes de conception d'API. Par exemple, il ne définit pas de formats de données à utiliser pour la réponse à renvoyer à l'utilisateur.ice<sup>27</sup>. La TEI n'est donc pas requise, peut-être de la critique faite par Smith et Blackwell des représentations hiérarchiques du texte. Si cela peut avoir un intérêt de ne pas vouloir « enfermer » le texte dans une hiérarchie qui lui est étrangère, ne pas utiliser la TEI pour une édition nativement numérique aujourd'hui est un choix pour le moins surprenant : ce standard dispose d'une très grande communauté d'utilisateur.ice.s, et est au cœur des humanités numériques centrées sur le texte. Le fait de ne pas imposer de format unique amène à beaucoup de variété dans les implémentations, ce qui s'oppose à un objectif de convergence des initiatives et ne garantit pas l'interopérabilité technique entre les API. De fait, les implémentations de CTS varient beaucoup, comme le montre l'exemple de *CapiTainS*, qui lui rend obligatoire l'utilisation du format XML dans les réponses. C'est pourquoi un nouveau standard a été développé : le DTS. Alors que le CTS existe depuis une décennie, la première version publique du CTS a été publiée en 2018. Elle est basée sur le même principe que le standard précédent : une sémantique standardisée permet à une URI de représenter un ou plusieurs (extraits) de texte. Plusieurs ajouts ont été faits à son prédécesseur. Le plus notable est l'obligation d'utiliser le XML-TEI<sup>28</sup>, ce qui garantit l'interopérabilité des ressources distribuées. Le standard cherche également à intégrer les meilleures pratiques en matière de développement Web, en utilisant des formats sémantiques (comme le JSON-LD) et en encourageant le développement d'API RESTful. Ce standard est également considérablement plus complexe que le CTS puisque l'API peut offrir plusieurs points d'accès aux textes ; ces points d'accès correspondent à différents niveaux de la collection (allant de la collection dans son ensemble au document). Bien que récent, le format a été implémenté par plusieurs projets, dont *TEI Publisher*, une application en ligne visant à publier des données de recherche en TEI.

Ces deux standards sont intéressants, mais ne sont pas adaptés à *KatAPI*. D'abord, leur implémentation complexifie considérablement la conception d'une API, et son utilisation est peut-être plus adaptée à des initiatives faisant converger plusieurs projets. Mais le principal problème posé par les trois standards présentés ici est qu'ils sont tous développés pour partager des textes réels, qui ont une existence « en dehors de l'API ». Or, *KatAPI* n'est pas seulement censée partager des textes issus du projet (c'est-à-dire,

---

27. B. Almas, T. Clérice, Hugh Cayless, V. Jolivet, Pietro Maria Liuzzo, Matteo Romanello, Jonathan Robie et Ian W. Scott, *Distributed Text Services (DTS) : a Community-built API to Publish and Consume Text Collections as Linked Data*, mars 2021, URL : <https://hal.archives-ouvertes.fr/hal-03183886> (visité le 01/01/2022), p. 2.

28. *Ibid.*, p. 3.



des catalogues ou des extraits de ceux-ci). Elle doit tout aussi bien partager des données de recherche issues du projet et construire des extraits du corpus en fonction des requêtes des utilisateur.ice.s. Cette API repose donc, comme nous le verrons, sur la création de documents *ex-nihilo* à partir de l'agrégation et de la transformation d'autres documents. Les critères bibliographiques servant à situer texte ou extrait dans une collection ne sont donc pas entièrement pertinents ici. L'API doit fonctionner comme un moteur de recherche plutôt que comme un catalogue dans lequel il est possible de sélectionner une ou plusieurs ressources. Cependant, les standards CTS et DTS présentent des bonnes pratiques, et de nombreux enseignements sont à en tirer. D'abord, il est nécessaire de proposer une réponse en XML-TEI, puisque ce standard de fait peut permettre de faire converger des données issues de différents projets. De CTS ressort l'intérêt de proposer plusieurs points d'accès au corpus, qui correspondent à différentes échelles dans celui-ci. Il n'existe pas une seule voie d'accès pour un texte, et c'est pourquoi ce principe sera suivi pour *KatAPI*. De CTS et DTS, il faut également retenir l'utilisation de l'URL comme d'un pointeur, à la sémantique clairement définie, vers les ressources pertinentes (un principe déjà énoncé dans le protocole HTTP). Enfin, des trois standards présentés ici, il faut retenir le principe de la *separation of concerns* (« séparation des préoccupations »). Ce principe central du développement applicatif est très bien exemplifié dans le fonctionnement d'OAI-PMH : la base de donnée est distincte de l'entrepôt OAI-PMH qui est lui-même distinct de la partie de l'API chargée de recevoir des requêtes et de transmettre des réponses au client. Ce principe permet de minimiser l'impact d'un processus sur un autre et de réduire les relations entre les parties d'un programme à un strict minimum. C'est un concept qui a été suivi dans l'architecture de *KatAPI*, où la réception et l'analyse des requêtes, la récupération des données, la construction des réponses et l'envoi de celle-ci à l'utilisateur.ice sont distinctes.

## 8.4 Pour qui sont ces API ? Qu'est-ce qui est demandé d'un tel outil ?

Cette dernière question n'est pas anodine. Il semblerait en effet que relativement peu d'API soient développées dans des contextes d'humanités numériques ; celles qui sont développées sont plutôt le fait de grandes institutions – comme le projet *Europeana*, ou la Bibliothèque nationale de France – et de plateformes partagées qui adoptent le standard OAI-PMH. Il existe également de nombreux services nécessitant l'utilisation d'API, comme les serveurs IIIF permettant le partage d'image, et les serveurs DTS qui commencent à se développer. La conception d'API semble être le fait d'infrastructures qui disposent de grands volumes de données et de ressources permettant de planifier et d'implémenter de tels outils sur le long terme.

Une étude réalisée dans le cadre du projet *Europeana* aide à situer le statut des API dans les humanités numériques. Cette étude offre un point de vue intéressant, puisqu'elle s'appuie sur des entretiens avec des chercheur.euse.s issue.e.s des humanités numériques et des humanités s'appuyant sur des méthodes quantitatives. Cette étude arrive à un constat pessimiste : les chercheur.euse.s ne sont pas des utilisateur.ice.s d'API<sup>29</sup>. Même lorsque leurs recherches s'appuient sur le traitement de données, voire sur l'usage de méthodes computationnelles (traitement statistique à l'aide de **python** ou **R**), ces personnes ont tendance à récupérer leurs données manuellement ; ce qui leur importe, c'est les données, et non la manière d'y accéder<sup>30</sup>. Les API sont, à l'inverse, beaucoup plus utilisées par des ingénieur.e.s et des personnes ayant eu une formation technique. Au delà de ce constat, l'étude propose des réflexions intéressantes sur le statut des API dans les humanités en général et numériques en particulier. Trois critères sont identifiés pour la décision d'utiliser une source de données en recherche<sup>31</sup> :

- Les données. Le critère de choix principal est la qualité et la complétude des données et des métadonnées.
- L'expertise technique nécessaire. Ce critère est relativement logique : un.e chercheur.euse n'utilisera une API que si il ou elle a le moyen de le faire. Un problème ici est la méconnaissance des API par les chercheur.euse.s, qui peuvent ne même pas connaître l'existence de tels outils. Cependant, l'étude remarque que, une fois qu'ils ont connaissance de telles sources, les chercheur.euse.s n'hésitent pas à se former à l'utilisation d'API.
- Les environnements sociaux et techniques des chercheur.euse.s. Le critère précédent décrit des capacités réelles, alors que celui-ci traite de la perception. Plus un.e chercheur.euse perçoit que des outils techniques sont difficiles d'accès, moins il ou elle est susceptible de les utiliser.

Pourquoi, alors, développer une API ? Au vu du volume du corpus, il peut être intéressant d'y avoir accès de façon automatisée. Avant le développement de cette API, la seule manière d'accéder à des données brutes était de les télécharger sur les dépôts GitHub du projet. Dans ce cas, il est uniquement possible de récupérer un fichier complet, c'est-à-dire un catalogue entier ou un jeu de données portant sur l'ensemble du corpus. La seule manière d'accéder à des données sélectionnées et filtrer en fonction de ses besoins est de passer par l'application Web « pour humain.e.s », mais celle-ci ne permet pas de télécharger les données brutes. Au sein du panorama d'outils développés pour le projet, les méthodes de diffusion de données manquaient. C'est pourquoi une API a été développée.

---

29. Jennifer Edmond et Vicky Garnett, « APIs and Researchers : The Emperor's New Clothes ? », *International Journal of Digital Curation*, 10-1 (21 mai 2015), p. 287-297, DOI : 10.2218/ijdc.v10i1.369, p. 288.

30. *Ibid.*, p. 290.

31. *Ibid.*, p. 292-294.

Plus largement, les API peuvent être particulièrement utiles pour un projet de recherche et pour l'application des principes FAIR en humanités numériques. Une API diminue, sur le long terme, le travail nécessaire pour diffuser des données. Dans le cas d'une API « réactive », qui construit des documents sur mesure, elle permet également d'explorer des sous-catégories, ou des sections du corpus. Une telle API peut permettre la création de nouveaux jeux de données à la demande, ce qui est d'un intérêt certain.

L'étude d'Edmond et Garnett ne doit pas être prise avec pessimisme, indiquant que l'outil que l'on développe ne sera jamais utilisé. Au contraire, les critères sociaux-techniques qu'elles mettent en avant derrière l'utilisation d'une API permettent de mieux adapter l'outil développé aux besoins. Il est nécessaire, lorsque cela est possible, de développer un outil qui diminue la distance perçue par les chercheur.euse.s. Pour cela il faut documenter l'application, de façon claire et, lorsque cela est possible, concise. Ensuite, il faut chercher à montrer que l'utilisation d'une API n'est pas très compliquée : il peut être intéressant de concevoir des tutoriels, ou comme cela a été fait pour *KatAPI*, des méthodes pour tester l'API en temps réel. Depuis la page de documentation, il est possible de composer un dictionnaire contenant les paramètres pour lancer une requête. Une fois la réponse reçue, le résultat s'affiche sur la page. De cette étude, il faut également retenir que une API ne sera utilisée que si elle offre l'accès à des données de qualité. Cela va dans le sens des principes FAIR : pour partager des données, celles-ci doivent être documentées et éditorialisées pour qu'un.e chercheur.euse sache quelles données seront reçues et envisage des manières de les utiliser.

Une autre étude<sup>32</sup>, menée sur des API généralistes, permet également de mieux comprendre le terrain pour de tels outils. Cette étude est basée sur l'analyse statistique des réponses obtenues pour 10955 *Business APIs*, terme qui recoupe des applications développées par des entreprises (comme l'API de *Twitter*), à but commercial (comme l'API de *PayPal*), mais aussi des API créées à but non-lucratifs (comme celle de *Wikipedia*). L'étude date d'il y a quelques années, et il est donc possible que les tendances de design d'API aient évoluées depuis. Les statistiques coïncident cependant avec mon expérience de ces outils. Les résultats obtenus permettent de mieux comprendre quelles sont les tendances actuelles en terme de design d'API. Les deux formats de réponse dominants sont le **JSON** et le **XML** : 81,08% des applications proposent l'un ou l'autre format. Seules 21,37% utilisent les deux formats, ce qui est le cas de *KatAPI*. Un autre fait intéressant est que une immense majorité d'API sont *RESTlike*<sup>33</sup>. Contrairement aux API *RESTful*, les API *RESTlike* sont inspirées par les principes REST sans nécessairement les appliquer

---

32. Andrea Janes, Tadas Remencius, Alberto Sillitti et Giancarlo Succi, « Towards Understanding of Structural Attributes of Web APIs Using Metrics Based on API Call Responses », dans *Open Source Software : Mobile Open Source Technologies*, dir. Luis Corral, Alberto Sillitti, Giancarlo Succi, Jelena Vlasenko et Anthony I. Wasserman, Series Title : IFIP Advances in Information and Communication Technology, Berlin, Heidelberg, 2014, t. 427, p. 83-92, DOI : 10.1007/978-3-642-55128-4\_11.

33. *Ibid.*, p. 3.

totalement. Ce terme s'est développé du fait de l'implémentation très inégale du standard. Le problème est que le terme *RESTlike* est nécessairement imprécis ; il est donc difficile de tirer une conclusion de cette statistique : dans une large mesure, les principes REST ont servi au développement de *KatAPI*. Quoi qu'il en soit, *KatAPI* semble s'inscrire, par les formats utilisés et les principes architecturaux suivis, dans le paysage actuel des API.

Les standards et études présentés ici permettent de mieux cerner le statut des API centrées sur le texte en humanités numériques ; elles permettent également d'identifier des tendances dans le design d'API ainsi que les besoins des utilisateur.ice.s. En ayant cette compréhension du paysage technique pour ces applications, il est possible de définir clairement le périmètre de l'API et la manière dont l'interaction client-serveur doit s'organiser, dans la sémantique des requêtes autant que dans les réponses obtenues.

# Chapitre 9

## Définir un périmètre : que partager, et comment partager ?

Après avoir présenté le paysage des API en général et le statut de ces outils dans les humanités numériques en particulier, ce chapitre décrit *KatAPI* du point de vue de l'utilisateur.ice : quelles données peuvent être obtenues via l'API, quels paramètres de recherche sont autorisés et quels sont les formats de données retournés par l'application. C'est donc les principes de la communication client-serveur de *KatAPI* qui sont ici présentés.

### 9.1 Grands principes pour l'architecture de l'API

Pour plus de clarté, ici sont présentés de façon concise l'ensemble des principes auxquels l'application doit se conformer ; ceux-ci forment une sorte de cahier des charges pour son implémentation technique.

- L'application supporte uniquement la méthode HTTP `GET` : il est uniquement possible de demander des données à l'application, et non d'ajouter des données à une base de données, par exemple.
- L'application supporte deux formats de réponse : `JSON` et `XML-TEI`. Le format par défaut est le `JSON`, pour sa légèreté et sa facilité de manipulation. Dans les deux cas, la structure des réponses est la même : un en-tête qui décrit le contexte de la réponse (à ne pas confondre avec les en-têtes HTTP) et un corps, qui contient les données récupérées.
- Trois niveaux d'accès aux données sont possibles. Il est possible de faire une requête pour un catalogue entier ; des statistiques sur un ou plusieurs catalogues ; une ou plusieurs entrées de catalogues.
- L'application suit le principe de séparation des préoccupations. L'interaction avec le client, le traitement des requêtes et la base de données sont des parties distinctes qui ne communiquent que lorsque cela est nécessaire.

- L’API suit les principes REST lorsque cela est possible. Elle est sans état et présente une interface uniforme. Cela veut dire que, pour les trois niveaux d’accès aux données, trois modèles de données seulement existent et sont adaptés en JSON et XML-TEI.
- L’application a un traitement strict des requêtes. Une requête n’est pas traitée si elle contient des paramètres non-autorisés (qui ne font pas partie de la sémantique définie pour l’interaction client-serveur), ou des valeurs associées à ces paramètres qui ne sont pas autorisées.
- L’application envoie des réponses précises et auto-descriptives, même en cas d’erreur. Plusieurs statuts HTTP sont définis pour différents cas de figures ; si erreur il y a, un message d’erreur complet et décrivant le contexte et les raisons de l’erreur est renvoyé par l’API. Cette erreur constitue un document JSON ou XML-TEI valide.

## 9.2 Quelles données partager ?

À différentes étapes, une grande variété de données ont été produites ; de plus, différentes manières d’accéder aux données ont été constituées par le projet. Par exemple, un système de réconciliation des différentes entrées de catalogue permet d’identifier lorsqu’un manuscrit est vendu plusieurs fois, et donc présent dans différents catalogues. Cela permet de regrouper ensemble différentes occurrences d’un même manuscrit. Ce type de données pourrait être partagé, de même que les données extraites de SPARQL et présentées dans la partie précédente. Cependant, la mise à disposition de différents types de données est difficile tout en se conformant au principe d’uniformité de l’interface du REST. En effet, chaque source de données a sa propre structure. Il donc est nécessaire, pour chaque source de données, de définir une représentation en JSON et en TEI. Pour se conformer à l’objectif originel du projet *MSS / Katabase* (étudier des catalogues en vente), il a été choisi de ne diffuser que des données issues des catalogues à travers l’API. Cela élimine donc les données issues de SPARQL. Pour ne se conserver qu’une seule représentation par type de données, il a été choisi de ne pas diffuser non plus de données sur les manuscrits « réconciliés », puisque cela amènerait à avoir deux modèles de données pour les manuscrits, selon qu’ils aient été réconciliés ou non. De plus, ce processus de réconciliation prend quelques secondes à s’exécuter, ce augmente le temps de réponse de l’API et peut devenir problématique si de nombreuses requêtes sont reçues en même temps par l’application.

Au final, trois niveaux d’accès aux corpus ont été définis. Ceux-ci correspondent à différents degrés de granularité dans le corpus. Chacun de ces formats, enfin, a un modèle de données qui lui est propre. Celui-ci est transposé à la fois en XML-TEI et en JSON.

- Premier niveau : catalogue complet (`cat_full`) : l’intégralité des données contenues dans un catalogue peut être requêtée. Dans ce cas, il n’est possible que de requêter

un catalogue à la fois. Celui-ci doit être clairement identifié à l'aide de son identifiant unique (`l@xml:id`). Il est renvoyé dans son intégralité au format XML uniquement : la traduction d'un fichier XML-TEI complet en JSON est complexe ; de plus, si un.e utilisateur.ice demande un catalogue de façon aussi ciblée, alors il ou elle souhaitera probablement un format plus complexe que le JSON et pourra tirer parti du balisage sémantique permise par la TEI.

- Deuxième niveau : statistiques sur une collection de catalogues (`cat_stat`) : des données statistiques sont retournées pour un ou plusieurs catalogues (prix moyen et médian des items, variance des prix au sein du catalogue, nombre d'items en vente...). En utilisant d'autres paramètres de recherche, il est possible de cibler précisément un sous-ensemble de catalogues, comme ceux issus de la *Revue des autographes*, ou ceux publiés dans une certaine tranche de dates. Ce degré de granularité permet donc de produire des sous-ensemble pour des études statistiques et thématiques ciblées.
- Troisième niveau : au niveau de l'entrée (`itm`). Ce degré permet d'accéder à une ou plusieurs entrées de catalogues, et donc de récupérer l'ensemble des informations contenues dans les catalogues pour ces entrées. Il est possible de cibler des entrées par nom auteur.ice, date d'écriture, date de vente ou par identifiant `@xml:id`, ce qui permet de ne retourner qu'un identifiant. Ce degré de granularité permet donc de constituer des jeux de données restreints, avec par exemple tous les manuscrits écrits par un.e auteur.ice à une certaine période.

Chacun de ces degrés d'accès aux données correspond à une source de données différente. Le premier niveau (`cat_full`) correspond à un document TEI représentant un catalogue dans son intégralité. Ce format permet donc l'accès aux sources directes du projet. Le deuxième niveau (`cat_stat`) donne l'accès à un JSON produit à partir d'une analyse statistique de tous les catalogues. Les entrées pertinentes de ce JSON sont traduites en XML-TEI si c'est le format de requête qui est demandée par l'utilisateur.ice. Enfin, le niveau `itm` correspond à deux sources de données : la première, c'est les catalogues eux-mêmes (au niveau de l'entrée, et non du catalogue complet) ; la deuxième, c'est une représentation en JSON des entrées de catalogues. Pour accélérer le fonctionnement de l'API lorsque la requête est faite au niveau de l'item, les données en JSON sont d'abord lues avant, si besoin, d'utiliser la TEI. Les trois degrés de granularité des données sont donc centraux à l'architecture de l'API : ils définissent trois bases de données à utiliser ; ils déterminent aussi trois représentation de données différentes (c'est-à-dire, trois structures pour les réponses de l'API). Grâce au principe de séparation des préoccupations, il est possible de rendre plus tard d'autres données accessibles et de permettre donc la diffusion de données depuis de nouvelles sources et la construction de nouvelles représentations. Les degrés de granularité déjà définis ne seront pas transformés par l'ajout d'autres degrés

de granularité pour l'accès au corpus. Puisque les données et leur représentation sont distinctes, il est également possible de modifier la structure des sources de données (la structure des catalogues, par exemple) sans que le format de sortie ne soit affecté. Il suffit pour ce faire de faire en sorte que la représentation des données soit la même malgré le changement de la source de données.

Couplés aux autres paramètres de recherche, ces trois degrés d'accès aux données permettent de construire des sous-ensembles du corpus de catalogues avec un grand degré de précision.

### **9.3 Codifier l'accès aux données : présentation des paramètres de recherche possibles**

Ici sont présentés les paramètres de recherche qui peuvent être utilisés pour récupérer des données, et comment ils permettent de filtrer ce qui est retourné par l'API.

### **9.4 Comment partager les données ? Principes suivis pour le partage d'informations, formats et structure des réponses de l'API**

Cette section décrit les formats de réponse de l'API, en JSON et XML-TEI, pour les différents paramètres de requête définis.

comment le header de la réponse permet de contenir des métadonnées tant voulues par le principe fair.



# Chapitre 10

## Implémentation et fonctionnement technique de *KatAPI*

Ce chapitre décrit le fonctionnement de l'API côté serveur : comment les données sont reçues, la manière dont elles sont traitées et comment les réponses sont construites (notamment la création automatisée de documents XML-TEI). Le système de gestion des erreurs est également présenté.

### 10.1 Présentation générale

Ici est présenté, schéma à l'appui, le fonctionnement technique de l'API et la chaîne de traitement depuis la réception des requêtes jusqu'à l'envoi d'une réponse.

### 10.2 Gestion des erreurs

L'API étant une application devant interagir avec des utilisateur.ice.s et gérer leurs requêtes, il n'est pas pensable qu'elle « plante » ou ne renvoie pas d'erreur. Tout un système de gestion des erreurs a donc été défini. Ces erreurs peuvent être classées en deux catégories : côté client (dus à des entrées invalides des utilisateur.ice.s) et côté serveur (erreurs inattendues qui empêchent l'exécution d'une requête). Si ces erreurs sont rencontrées, des réponses sont construites en JSON ou XML-TEI (selon le format demandé par l'utilisateur.ice.s). Celles-ci décrivent l'erreur et permettent donc aux utilisateur.ice.s de corriger les erreurs qui ont pu avoir lieu côté client.

### 10.3 Garantir le bon fonctionnement de l'application

De la même manière que les systèmes de gestion des erreurs garantissent que l'application continue à fonctionner même en cas d'erreur, les tests garantissent que l'API soit

fonctionnelle, même dans « des conditions extrêmes ». Les réponses, formats et données retournées sont donc testés. C'est ce protocole de test qui est présenté ici.

# Conclusion

Ainsi, en s'appuyant sur une connaissance du corpus, il est possible d'approximer une véritable compréhension « humaine » d'un texte semi structuré, en utilisant uniquement des expressions régulières. Si il n'est pas possible d'identifier le sens des différents éléments, il est possible de localiser et d'extraire les éléments signifiants. Il apparaît alors que ces éléments signifiants sont, d'une certaine manière, des formes, ou des motifs qui peuvent être détectés automatiquement ; la signification peut être identifiée en fonction du type de donnée extraite (un nombre est une date). Mais, dans un corpus semi-structuré, il surtout est possible d'inférer du « sens » d'un élément à partir de sa position dans le texte. C'est dans ce cadre qu'un encodage en XML-TEI des documents prend tout son intérêt : les différentes parties d'un texte sont balisées sémantiquement, ce qui est essentiel au traitement automatisé du texte : il devient possible de traiter le texte non pas dans son intégralité, mais au niveau d'une unité sémantique très précise, comme le nom donné à un manuscrit (le `tei:name`) ou une brève description (le `tei:trait`)

Le fait de travailler sur des corpus semi-structurés permet également de mettre en approche des méthodes alternatives, « low-tech », pour le traitement du langage et l'analyse textuelle dans des projets de recherche. Avec une bonne connaissance du corpus et des données structurées, il est possible de se baser entièrement sur des solutions techniquement « simples », telles que la détection de motifs à l'aide d'expressions régulières. Plutôt que de s'appuyer sur de l'apprentissage machine, et donc des solutions complexes, il est intéressant de s'appuyer sur des méthodes alternatives et de développer des algorithmes plus complexes, basés sur une approche modulaire, qui sont capables de cibler très précisément les éléments d'un texte pertinent pour des problématiques de recherche précises. Ces méthodes de basse technologie, qui s'intéressent à un usage intelligent des technologies plutôt qu'à des technologies « intelligentes » forment un contre modèle qui gagnerait à être mis en avant dans les projets de recherche. L'intelligence artificielle devenue un outil central et pertinent dans les humanités numériques – à la base de la reconnaissance optique de caractères. Mais elle ne doit pas seulement être une nouveauté que les projets de recherche doivent à tout prix incorporer. C'est également une technologie extrêmement polluante<sup>1</sup>,

---

1. Strubell, Ganesh et McCallum (E. Strubell, A. Ganesh et A. McCallum, « Energy and Policy Considerations for Deep Learning in NLP »...) ont étudié l'impact écologique des modèles d'apprentissage profond en traitement du langage naturel. En étudiant les ressources nécessaires à l'entraînement d'un modèle, les auteurs montrent (*Ibid.*, p. 1) que la mise au point d'un modèle *Transformer* peut émettre

utilisée de façon massive par des grands groupes industriels<sup>2</sup> et un argument de vente pour de nombreuses start-up, alors que les applications pratiques de l'intelligence artificielle soulèvent de graves problèmes éthiques<sup>3</sup>. La recherche en humanités numériques, fondée sur une connaissance des techniques et des documents utilisés, devrait fournir un espace pour remettre en question de telles technologies et pour développer des solutions alternatives, où la connaissance scientifique supplante l'intelligence machine. De façon modeste, le projet *MSS / Katabase* offre une telle remise en question : en s'appuyant sur une bonne connaissance de corpus textuels semi-structurés, le projet a développé une chaîne de traitement entièrement basée sur la détection de motifs. S'appuyer sur des méthodes de basse technologie ne vaut bien sûr pas seulement comme une déclaration de bonne foi : plus une solution est techniquement simple et accessible, plus elle est réutilisable, adaptable à d'autres jeux de données et plus facilement elle peut évoluer pour incorporer de nouvelles fonctionnalités.

pollution des large language models : <https://arxiv.org/abs/1906.02243>

— parler du minimal computing ??

---

jusqu'à 626155 livres de CO<sub>2</sub>, soit cinq fois la consommation d'une voiture durant toute sa vie (126000 livres de CO<sub>2</sub>) et 17 fois les émissions de CO<sub>2</sub> émises par une personne américaine en un an (36156 livres). Il.elle.s rappellent également que les modèles d'apprentissage machine les plus performants sont aussi les plus gourmands en ressources et qu'un modèle doit être entraîné parfois plusieurs milliers de fois pour atteindre des scores satisfaisants (*Ibid.*). Avec l'augmentation générale des performances des modèles de traitement du langage naturel, l'augmentation marginale de leur performance est de plus en plus coûteuse en ressources : l'augmentation du score BLEU d'un modèle de traduction automatique de 0.1 points de performance (pour atteindre un score de 29,7) a coûté 150000 dollars en ressources environnementales et en coût de calcul (*Ibid.*, p. 4).

2. *Google*, et sa filiale *Youtube*, s'appuient depuis le milieu des années 2010 de plus en plus massivement sur l'intelligence artificielle (P. Covington, J. Adams et E. Sargin, « Deep Neural Networks for YouTube Recommendations »..., p. 1). L'algorithme de suggestion de vidéos sur *Youtube* repose entièrement sur de l'apprentissage machine (*ibid.*) : en fonction de l'historique de l'utilisateur.ice, l'algorithme « apprend » à suggérer du contenu pertinent. L'usage massif de telles technologies est particulièrement polluant : pour pouvoir suggérer des vidéos pertinentes, il faut, à chaque connexion de l'utilisateur.ice, consulter deux très grandes bases de données : dans un premier temps, l'algorithme sélectionne toutes les vidéos potentiellement pertinentes, depuis l'ensemble des vidéos disponibles sur le site ; ensuite, ces vidéos sont classées en consultant une seconde base de données et en filtrant les résultats obtenus par le comportement récent de l'utilisateur.ice (ce qu'il ou elle a recherché, combien de temps il ou elle a passé à regarder une vidéo, sa réaction à celle-ci...). Le 2 décembre 2020, *Google* a licencié la docteure Timnit Gebru, alors à la tête de l'équipe spécialisée en éthique de l'intelligence artificielle, suite à la rédaction d'un article sur les impacts écologiques de l'utilisation massive par l'entreprise de modèles de traitement automatisé du langage. Ce phénomène alarmant rappelle au passage la non-neutralité des publications scientifiques provenant de cette entreprise, dont deux ont été citées dans le présent mémoire (T. Mikolov, K. Chen, G. Corrado, *et al.*, « Efficient Estimation of Word Representations in Vector Space »... ; P. Covington, J. Adams et E. Sargin, « Deep Neural Networks for YouTube Recommendations »...).

3. À ce sujet, voir par exemple le projet *Awful AI* de David Dao, qui liste une quantité impressionnante d'utilisations préoccupantes de telles technologies.

# Bibliographie

## Projet *MSS* / *Katabase*

- CORBIÈRES (Caroline), *Du catalogue au fichier TEI. Création d'un workflow pour encoder automatiquement en XML-TEI des catalogues d'exposition*. Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2020, URL : [https://github.com/carolinecorbieres/Memoire\\_TNAH](https://github.com/carolinecorbieres/Memoire_TNAH) (visité le 13/06/2022).
- GABAY (Simon), RONDEAU DU NOYER (Lucie) et KHEMAKHEM (Mohamed), « Selling autograph manuscripts in 19th c. Paris : digitising the Revue des Autographes », dans *IX Convegno AIUCD*, Milan, 2020, URL : <https://hal.archives-ouvertes.fr/hal-02388407> (visité le 13/06/2022).
- GABAY (Simon), RONDEAU DU NOYER (Lucie), GILLE LEVENSON (Matthias), PETKOVIC (Ljudmila) et BARTZ (Alexandre), « Quantifying the Unknown : How many manuscripts of the marquise de Sévigné still exist ? », dans *Digital Humanities DH2020*, Ottawa, 2020, URL : <https://hal.archives-ouvertes.fr/hal-02898929/> (visité le 13/06/2022).
- JANÈS (Juliette), *Du catalogue papier au numérique. Une chaîne de traitement ouverte pour l'extraction d'information issue de documents structurés*, Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2021, URL : [https://github.com/Juliettejns/Memoire\\_TNAH](https://github.com/Juliettejns/Memoire_TNAH) (visité le 13/06/2022).
- KHEMAKHEM (Mohamed), ROMARY (Laurent), GABAY (Simon), BOHBOT (Hervé), FRONTINI (Francesca) et LUXARDO (Giancarlo), « Automatically Encoding Encyclopedic-like Resources in TEI », dans *The annual TEI Conference and Members Meeting*, Tokyo, 2018, URL : <https://hal.archives-ouvertes.fr/hal-01819505> (visité le 13/06/2022).
- « Information Extraction Workflow for Digitised Entry-based Documents. » Dans *DARIAH Annual event 2020*, Zagreb, 2020, URL : <https://hal.archives-ouvertes.fr/hal-02508549> (visité le 13/06/1997).
- RONDEAU DU NOYER (Lucie), *Encoder automatiquement des catalogues en XML-TEI. Principes, évaluation et application à la revue des autographes de la librairie Cha-*

ravay, Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2019, URL : <https://github.com/lairaines/M2TNAH> (visité le 13/06/2022).

RONDEAU DU NOYER (Lucie), GABAY (Simon), KHEMAKHEM (Mohamed) et ROMARY (Laurent), « Scaling up Automatic Structuring of Manuscript Sales Catalogues », dans *TEI 2019 : What is text, really? TEI and beyond*, Graz, 2019, URL : <https://hal.archives-ouvertes.fr/hal-02272962> (visité le 13/06/2022).

## Édition numérique, traitement automatisé et analyse de texte

AGIRRE (Eneko), BARRENA (Ander), LACALLE (Oier Lopez de), SOROA (Aitor), FERNANDO (Samuel) et STEVENSON (Mark), « Matching Cultural Heritage items to Wikipedia », dans *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, dir. Nicoletta Calzolari (Conference Chair), *et al.*, Istanbul, Turkey, 2012.

BLEEKER (Elli), HAENTJENS DEKKER (Ronald) et BUITENDIJK (Bram), « Texts as Hypergraphs : An Intuitive Representation of Interpretations of Text », *Journal of the Text Encoding Initiative*, 14 (2021), DOI : <https://doi.org/10.4000/jtei.3919>.

BOEGLIN (Noémie), DEPEYRE (Michel), JOLIVEAU (Thierry) et LE LAY (Yves-François), « Pour une cartographie romanesque de Paris au XIXe siècle. Proposition méthodologique », dans *Conférence Spatial Analysis and GEomatics*, Nice, France, 2016 (Actes de la conférence SAGEO'2016 - Spatial Analysis and GEomatics), URL : <https://hal.archives-ouvertes.fr/hal-01619600>.

BRANDO (Carmen), ABADIE (Nathalie) et FRONTINI (Fontini), « Évaluation de la qualité des sources du web de données pour la résolution d'entités nommées », *Ingénierie des systèmes d'information*, 21-5 (28 déc. 2016), p. 31-54, DOI : 10.3166/isi.21.5-6.31-54.

BRANDO (Carmen), FRONTINI (Francesca) et GANASCIA (Jean-Gabriel), « REDEN : Named Entity Linking in Digital Literary Editions Using Linked Data Sets », *Complex Systems Informatics and Modeling Quarterly*-7 (29 juil. 2016), p. 60-80, DOI : 10.7250/csimq.2016-7.04.

BURNARD (Lou), « What is TEI conformance, and why should you care? », *Journal of the Text Encoding Initiative*, 12 (2019), DOI : <https://doi.org/10.4000/jtei.1777>.

BURNARD (Lou), SCHÖCH (Christian) et ODEBRECHT (Carolyn), « In search of comity : TEI for distant reading », *Journal of the Text Encoding Initiative*, 14 (2021), DOI : <https://doi.org/10.4000/jtei.3500>.

- CHAGUÉ (Alix) et CLÉRICE (Thibault), « Sharing HTR datasets with standardized metadata : the HTR-United initiative », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris, France, 2022, URL : <https://hal.inria.fr/hal-03703989>.
- CHRISTENSEN (Kelly), GABAY (Simon), PINCHE (Ariane) et CAMPS (Jean-Baptiste), « SegmOnto – A Controlled Vocabulary to Describe Historical Textual Sources », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris : École nationale des Chartes, 2022.
- Expression régulière*, Wikipedia. L’encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Expression\\_r%C3%A9guli%C3%A8re](https://fr.wikipedia.org/wiki/Expression_r%C3%A9guli%C3%A8re) (visité le 10/07/2022).
- FRONTINI (Francesca), BRANDO (Carmen), RIGUET (Marine), JACQUOT (Clémence) et JOLIVET (Vincent), « Annotation of Toponyms in TEI Digital Literary Editions and Linking to the Web of Data », *Matlit Revista do Programa de Doutorado em Materialidades da Literatura*, 4-2 (11 juil. 2016), p. 49-75, DOI : 10.14195/2182-8830\_4-2\_3.
- GABAY (Simon), CAMPS (Jean-Baptiste), PINCHE (Ariane) et JAHAN (Claire), « SegmOnto : common vocabulary and practices for analysing the layout of manuscripts (and more) », dans *1st International Workshop on Computational Paleography (IWCP@ICDAR 2021)*, Lausanne, 2021, URL : <https://hal.archives-ouvertes.fr/hal-03336528> (visité le 29/06/2022).
- MENDES (Pablo N.), JAKOB (Max), GARCÍA-SILVA (Andrés) et BIZER (Christian), « DBpedia spotlight : shedding light on the web of documents », dans *Proceedings of the 7th International Conference on Semantic Systems - I-Semantics '11*, Graz, Austria, 2011, p. 1-8, DOI : 10.1145/2063518.2063519.
- RENEAR (Allen), MYLONAS (Elli) et DURAND (David), « Refining our Notion of What Text Really Is : The Problem of Overlapping Hierarchies », dans *Research in Humanities Computing*, dir. Nancy Ide et Susan Hockey, Oxford, 1996, URL : <https://cds.library.brown.edu/resources/stg/monographs/ohco.html> (visité le 22/08/2022).
- RUIZ FABO (Pablo) et POIBEAU (Thierry), « Mapping the Bentham Corpus : Concept-based Navigation », *Journal of Data Mining & Digital Humanities*, Atelier Digit\_Hum (Data deluge : which skills for...[ 2019]), p. 5044, DOI : 10.46298/jdmdh.5044.
- SAHLE (Patrick), « Digital modelling. Modelling the digital edition », dans *Medieval and modern manuscript studies in the digital age*, London/Cambridge, 2016, URL : [https://dixit.uni-koeln.de/wp-content/uploads/2015/04/Camp1-Patrick\\_Sahle\\_-\\_Digital\\_Modelling.pdf](https://dixit.uni-koeln.de/wp-content/uploads/2015/04/Camp1-Patrick_Sahle_-_Digital_Modelling.pdf).

- SAHLE (Patrick), « What is a scholarly digital edition ? », dans *Digital scholarly editing. Theories and practices*, dir. Matthew James Driscoll et Elena Pierazzo, Cambridge, 2016, URL : <http://books.openedition.org/obp/3397> (visité le 10/07/2022).
- SAYERS (Jentery), *Minimal definitions*, Minimal Computing - a working group of GO ::DH, 2 oct. 2016, URL : <http://go-dh.github.io/mincomp/thoughts/2016/10/02/minimal-definitions/> (visité le 10/06/2022).
- SCHWEYER (Anne-Valérie), BURIE (Jean-Christophe) et NGUYEN (Tien Nam), « Analyse, Reconnaissance et Indexation des manuscrits CHAM », dans *Documents anciens et reconnaissance automatique des écritures manuscrites*, Paris, 2022.
- SOUDANI (Aicha), MEHERZI (Yosra), BOUHAFS (Asma), FRONTINI (Francesca), BRANDO (Carmen), DUPONT (Yoann) et MÉLANIE-BECQUET (Frédérique), « Adaptation et évaluation de systèmes de reconnaissance et de résolution des entités nommées pour le cas de textes littéraires français du 19ème siècle », dans *Atelier Humanités Numériques Spatialisées (HumaNS'2018)*, Montpellier, 2018, URL : <https://hal.archives-ouvertes.fr/hal-01925816> (visité le 10/08/2022).
- STOKES (Peter A.), KIESSLING (Benjamin), STÖKL BEN EZRA (Daniel), TISSOT (Robin) et GARGEM (Hassane), « The eScriptorium VRE for Manuscript Cultures », *Classics@ Journal*, 18 (2021), URL : <https://classics-at.chs.harvard.edu/classics18-stokes-kiessling-stokl-ben-ezra-tissot-gargem/> (visité le 14/07/2022).
- STRUBELL (Emma), GANESH (Ananya) et MCCALLUM (Andrew), « Energy and Policy Considerations for Deep Learning in NLP », *CoRR*, abs/1906.02243 (2019), DOI : <https://doi.org/10.48550/arXiv.1906.02243>.
- TEI CONSORTIUM, *P5 : Guidelines for Electronic Text Encoding and Interchange*, Text Encoding Initiative, Version 4.4.0, 2022, URL : <https://tei-c.org/release/doc/tei-p5-doc/en/html/index.html>.

## Technologies du web : généralités

- BOECKHOUT (Martin), ZIELHUIS (Gerhard A.) et BREDENOORD (Annelien L.), « The FAIR guiding principles for data stewardship : fair enough ? », *European Journal of Human Genetics*, 26–7 (juil. 2018), p. 931-936, DOI : 10.1038/s41431-018-0160-0.
- COVINGTON (Paul), ADAMS (Jay) et SARGIN (Erme), « Deep Neural Networks for YouTube Recommendations », dans *Proceedings of the 10th ACM Conference on Recommender Systems*, New York, 2016, URL : <https://research.google/pubs/pub45530/> (visité le 10/06/2022).
- MIKOLOV (Tomas), CHEN (Kai), CORRADO (Greg) et DEAN (Jeffrey), « Efficient Estimation of Word Representations in Vector Space » (, 2013), Publisher : arXiv Version Number : 3, DOI : 10.48550/ARXIV.1301.3781.



*Moteur de recherche*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Moteur\\_de\\_recherche](https://fr.wikipedia.org/wiki/Moteur_de_recherche) (visité le 17/07/2022).

*Namespace*, Wikipedia. L'encyclopédie libre, 2022, URL : <https://en.wikipedia.org/wiki/Namespace> (visité le 05/05/2022).

*Wikimedia Foundation\$Technology*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://en.wikipedia.org/wiki/Wikimedia\\_Foundation#Technology](https://en.wikipedia.org/wiki/Wikimedia_Foundation#Technology) (visité le 01/08/2022).

## Technologies du web : Web sémantique

BECKETT (Dave), BROEKSTRA (Jeen) et HAWKE (Sandro), *SPARQL Query Results XML Format (Second Edition)*. W3C Recommendation 21 March 2013, W3C, 21 mars 2013, URL : <https://www.w3.org/TR/2013/REC-rdf-sparql-XMLres-20130321/> (visité le 10/05/2022).

ERXLEBEN (Fredo), GÜNTHER (Michael), KRÖTZSCH (Markus), MENDEZ (Julian) et VRANDEČIĆ (Denny), « Introducing Wikidata to the Linked Data Web », dans *The Semantic Web – ISWC 2014*, dir. Peter Mika, *et al.*, Series Title : Lecture Notes in Computer Science, Cham, 2014, t. 8796, p. 50-65, DOI : 10.1007/978-3-319-11964-9\_4.

*Graph database*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://en.wikipedia.org/wiki/Graph\\_database](https://en.wikipedia.org/wiki/Graph_database) (visité le 12/06/2022).

HARRIS (Steve), SEABORNE (Andy) et PRUD'HOMMEAUX (Eric), *SPARQL 1.1 Query Language*. W3C Recommendation 21 March 2013, W3C, 2013, URL : <http://www.w3.org/TR/2013/REC-sparql11-query-20130321/> (visité le 10/06/2022).

MÜLLER-BIRN (Claudia), KARRAN (Benjamin), LEHMANN (Janette) et LUCZAK-RÖSCH (Markus), « Peer-production system or collaborative ontology engineering effort : what is Wikidata ? », dans *Proceedings of the 11th International Symposium on Open Collaboration*, San Francisco California, 2015, p. 1-10, DOI : 10.1145/2788993.2789836.

*Ontology (information science)*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://en.wikipedia.org/wiki/Ontology\\_\(information\\_science\)](https://en.wikipedia.org/wiki/Ontology_(information_science)) (visité le 01/02/2022).

*Query optimization*, Wikidata, 2022, URL : [https://www.wikidata.org/wiki/Wikidata:SPARQL\\_query\\_service/query\\_optimization](https://www.wikidata.org/wiki/Wikidata:SPARQL_query_service/query_optimization) (visité le 06/06/2022).

## Technologies du web : API

ALMAS (Bridget) et CLÉRICE (Thibault), « Continuous Integration and Unit Testing of Digital Editions », *Digital Humanities Quarterly*, 11-4 (févr. 2018), Publisher :

- Alliance of Digital Humanities, URL : <https://hal.archives-ouvertes.fr/hal-01709868> (visité le 12/07/2022).
- ALMAS (Bridget), CLÉRICE (Thibault), CAYLESS (Hugh), JOLIVET (Vincent), LIUZZO (Pietro Maria), ROMANELLO (Matteo), ROBIE (Jonathan) et SCOTT (Ian W.), *Distributed Text Services (DTS) : a Community-built API to Publish and Consume Text Collections as Linked Data*, mars 2021, URL : <https://hal.archives-ouvertes.fr/hal-03183886> (visité le 01/01/2022).
- API :Main page, MediaWiki, 2022, URL : [https://www.mediawiki.org/wiki/API:Main\\_page](https://www.mediawiki.org/wiki/API:Main_page) (visité le 10/06/2022).
- CLÉRICE (Thibault), ALMAS (Bridget), CAYLESS (Hugh), JOLIVET (Vincent), MORLOCK (Emmanuelle), ROBIE (Jonathan), TAUBER (James), WITT (Jeffrey) et LIUZZO (Pietro), « From File Interoperability to Service Interoperability : The Distributed Text Services », dans *TEI 2018*, Tokyo, Japan, 2018, URL : <https://hal.archives-ouvertes.fr/hal-02196659> (visité le 09/08/2022).
- EDMOND (Jennifer) et GARNETT (Vicky), « APIs and Researchers : The Emperor's New Clothes ? », *International Journal of Digital Curation*, 10–1 (21 mai 2015), p. 287-297, DOI : 10.2218/ijdc.v10i1.369.
- FIELDING (Roy), *Architectural Styles and the Design of Network-based Software Architectures*, Version : PDF edition, 1-column for viewing online, Thèse de doctorat, Irvine, University of California, 2000, URL : <https://www.ics.uci.edu/~fielding/pubs/dissertation/top.htm> (visité le 05/08/2022).
- GOYET (Samuel), *De briques et de blocs. La fonction éditoriale des interfaces de programmation (API) web : entre science combinatoire et industrie du texte*. Theses, Paris IV Sorbonne, 2017, URL : <https://tel.archives-ouvertes.fr/tel-01665406> (visité le 10/08/2022).
- Interface de programmation, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Interface\\_de\\_programmation](https://fr.wikipedia.org/wiki/Interface_de_programmation) (visité le 27/07/2022).
- JANES (Andrea), REMENCIUS (Tadas), SILLITTI (Alberto) et SUCCI (Giancarlo), « Towards Understanding of Structural Attributes of Web APIs Using Metrics Based on API Call Responses », dans *Open Source Software : Mobile Open Source Technologies*, dir. Luis Corral, *et al.*, Series Title : IFIP Advances in Information and Communication Technology, Berlin, Heidelberg, 2014, t. 427, p. 83-92, DOI : 10.1007/978-3-642-55128-4\_11.
- MURDOCK (Jaimie) et ALLEN (Colin), « InPhO for All : Why APIs Matter », *Journal of the Chicago Colloquium on Digital Humanities and Computer Scienc*, 1–3 (2011), URL : <http://www.jamram.net/docs/jdhcs11-paper.pdf> (visité le 12/07/2022).
- PRIME-CLAVERIE (Camille) et MAHÉ (Annaïg), « Le défi de l'interopérabilité entre plates-formes pour la construction de savoirs augmentés en sciences humaines et

- sociales », dans *Ecrilecture augmentée dans les communautés scientifiques*, 2017, URL : [https://archivesic.ccsd.cnrs.fr/sic\\_01511618](https://archivesic.ccsd.cnrs.fr/sic_01511618) (visité le 29/08/2022).
- SMITH (David Neel) et BLACKWELL (Christopher W.), « Four URLs, limitless apps : Separation of concerns in the Homer Multitext architecture », dans *Donum natalicium digitaliter confectum Gregorio Nagy septuagenario a discipulis collegis familiaribus oblatum : A Virtual Birthday Gift Presented to Gregory Nagy on Turning Seventy by His Students, Colleagues, and Friends*, 2012, URL : <https://chs.harvard.edu/d-n-smith-c-w-blackwell-four-urls-limitless-apps-separation-of-concerns-in-the-homer-multitext-architecture/> (visité le 22/08/2022).
- WANG (Sean), BELOUIN (Pascal) et CHEN (Shih-Pei), « Research Infrastructure for the Study of Eurasia (RISE) : Towards a Flexible and Distributed Digital Infrastructure for Resource Access via Standardized Apis and Metadata », dans *DADH 2018. 9th International Conference of Digital Archives and Digital Humanities*, Taipei, 2018, p. 21-37, URL : [https://pure.mpg.de/pubman/faces/ViewItemOverviewPage.jsp?itemId=item\\_3033461](https://pure.mpg.de/pubman/faces/ViewItemOverviewPage.jsp?itemId=item_3033461) (visité le 01/07/2022).

## Visualisation et design d'interfaces

- AGAMBEN (Giorgio), *Qu'est-ce qu'un dispositif ?*, Stanford, 2006 (Rivages Poche / Petite Bibliothèque).
- ALBOUY (Ségolène), *Médiation des données de la recherche. Élaboration d'une plateforme en ligne pour une base de tables astronomiques anciennes*. Mémoire pour le diplôme de master "Technologies numériques appliquées à l'histoire", Paris, École nationale des Chartes, 2019, URL : <https://github.com/Segolene-Albouy/Memoire-TNAH2019> (visité le 13/06/2022).
- DRUCKER (Johanna), *Visualisation : l'interprétation modélisante*, Paris, 2020 (Esthétique des données, 03).
- HENTSCHEL (Klaus), *Visual cultures in science and technology. A Comparative History*, Oxford, 2014.
- MASURE (Anthony), *Design et humanités numériques*, Paris, 2017 (Esthétique des données, 01).
- ORAIN (Olivier), « Le Laboratoire de cartographie dans le contexte de développement des sciences sociales et humaines, des années 1950 aux années 1970 », dans *Design graphique, recherche et patrimoine des sciences sociales. Le Laboratoire de graphique de Jacques Bertin*, Pierrefitte-sur-Seine : Archives nationales, 2021, URL : <https://dai.ly/x85jbir> (visité le 14/07/2022).
- RENON (Anne-Lyse), « "Design graphique" et "objectivité", la question des méta-altas », dans *Voir l'architecture. Contribution du design à la construction des savoirs*, dir. Annick Lantenois et Gilles Rouffineau, Paris, Grenoble, 2015, p. 71-81.

RENON (Anne-Lyse), *Design et esthétique dans les pratiques de la science*, Thèse de doctorat, Paris, École des Hautes Études en Sciences Sociales, Institut Marcel Mauss, 2016, URL : [https://www.academia.edu/36754513/Design\\_et\\_esth%C3%A9tique\\_dans\\_les\\_pratiques\\_de\\_la\\_science](https://www.academia.edu/36754513/Design_et_esth%C3%A9tique_dans_les_pratiques_de_la_science) (visité le 08/06/2022).

## Marché de l'art, économétrie et statistiques

*F-score*, Wikipedia. L'encyclopédie libre, 2022, URL : <https://en.wikipedia.org/wiki/F-score> (visité le 12/06/2022).

MAUPEOU (Félicie Faizand de) et SAINT-RAYMOND (Léa), « Les “marchands de tableaux” dans le Bottin du commerce : une approche globale du marché de l'art à Paris entre 1815 et 1955 », *Artl@s Bulletin*, 2 (2013), URL : <https://hal.archives-ouvertes.fr/hal-02986371> (visité le 20/06/1997).

PIKETTY (Thomas), *Les hauts revenus en France au XXe siècle : inégalités et redistributions, 1901-1998*, Paris, 2001.

*Précision et rappel*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Pr%C3%A9cision\\_et\\_rappel](https://fr.wikipedia.org/wiki/Pr%C3%A9cision_et_rappel) (visité le 13/06/2022).

*Régression linéaire*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/R%C3%A9gression\\_lin%C3%A9aire](https://fr.wikipedia.org/wiki/R%C3%A9gression_lin%C3%A9aire) (visité le 10/07/2022).

SAINT-RAYMOND (Léa), « Revisiting Harrison and Cynthia White's Academic vs. Dealer-Critic System », *Arts*, 8–3 (2019), DOI : <https://doi.org/10.3390/arts8030096>.

# Glossaire

**Open source** L'*open source*, ou « code ouvert » définit à la fois un statut légal et un mouvement social. L'*open source* correspond à un logiciel ou programme dont le code est librement consultable en ligne et réutilisable selon différentes licences plus ou moins permissives, comme la *GNU GPL v3.0* ou la *MIT License*. Ce statut légal a beaucoup gagné en popularité, en partie pour des raisons pragmatiques (un code ouvert est maintenu gratuitement, ce qui diminue les coûts d'entreprises) et symboliques (les GAFAM se sont beaucoup associées à ce courant en raison de l'image positive qu'il donne). Le statut est également associé à plusieurs mouvements en faveur du partage et de la réutilisation des connaissances, tel que la science ouverte. Certaines franges de l'*open source* sont plus radicales et ont une conception ouvertement politique de l'accessibilité au code, accessibilité qui permettrait aux utilisateur.ice.s de regagner en autonomie, de développer de nouveaux modèles communautaires et de reprendre le contrôle sur leurs comportements en ligne et leur consommation médiatique. Parmi les représentant.e.s les plus proéminent.e.s de ce courant se trouve Richard Stallman, créateur de la *GNU Foundation* et du système d'exploitation GNU, à la base de la majorité des systèmes d'exploitation utilisant Linux.. 69, 75

**API** Une interface de programmation d'application (*Application Programming Interface* en anglais) est un protocole qui permet à un programme de communiquer avec un autre programme. Ce protocole documenté, correspond à un ensemble d'opérations permettant à un programme « consommateur » d'utiliser des fonctionnalités d'un programme « fournisseur », comme par exemple de récupérer ou d'envoyer des données au fournisseur.. 20

**Base de données en graphe** Une base de données en graphe sert à stocker des données reliées entre elles au sein d'un graphe ou d'un réseau. Les données correspondent aux nœuds du graphe et les relations sont représentées sous forme d'arrêtes. Souvent, données et relations disposent d'identifiants uniques définis grâce à des ontologies. L'interaction avec de telles bases de données est permise grâce à des langages de requêtes, comme SPARQL. La représentation en graphe a pour but de mettre l'accent sur l'interrelation entre les données et de permettre des liens complexes entre elles.<sup>4</sup>.

---

4. *Graph database*, Wikipedia. L'encyclopédie libre, 2022, URL : <https://en.wikipedia.org/>

81, 82, 87

**CTS** *Canonical Text Services* (CTS) est une API et une spécification pour identifier et diffuser des citations au sein de textes encodés. L'API spécifie la manière d'identifier un texte, un témoin, un passage ou encore une série de textes à l'aide d'un URN (un identifiant unique analogue aux URI). Le CTS est une API, mais ne spécifie pas comment une API CTS doit être implémentée : ce qui importe, c'est que l'API CTS permette d'accéder à des extraits de textes grâce à leurs citations<sup>5</sup>. 107

**Dictionnaire** Un dictionnaire est un format de données structuré en python qui associe à une donnée – dite clé – une ou plusieurs données nommées valeurs.. 25, 37, 40, 44, 47, 48, 50, 54, 63, 70–72, 119

**DTS** Comme le CTS, le *Distributed Text Services* (DTS) est une API visant à identifier et permettre l'identification et la diffusion de texte dans une architecture client-serveur. Cette API, développée à partir de 2018 sur un modèle communautaire, demande une réponse en XML-TEI et des métadonnées encodées en JSON-LD ; ces spécifications techniques la distinguent de CTS, qui ne repose pas sur un format particulier. Comme le CTS, DTS ne définit pas comment implémenter l'API<sup>6</sup>. 107

**Expression régulière** Une expression régulière, ou expression rationnelle est une chaîne de caractère, écrite selon une syntaxe précise, qui permet de détecter des motifs dans du texte<sup>7</sup>. Les expressions régulières s'appuient sur la classification des caractères en classes (minuscules, majuscules, chiffres, espaces), sur des structures alternatives (un caractère ou un autre) et exclusives (un caractère n'ayant pas certaines propriétés) pour repérer des motifs. Par exemple, « 2022 » correspond au motif : `\d4`, soit « quatre chiffres à la suite ». Une adresse mail est également un motif : `[^(@|\s)]+@[^(@|\s)]+`, soit « plusieurs caractères qui ne sont ni un espace une arabase, suivi d'une arabase, suivi de plusieurs caractères qui n'est ni un espace ni une arabase ».. 24, 25, 35, 39, 50, 54, 59, 60, 127

**FAIR** Les principes *Findable Accessible Interoperable Reusable* (FAIR) correspondent à un ensemble de préconisations simples pour le partage de données dans un contexte de sciences ouvertes<sup>8</sup>. 107

---

wiki/Graph\_database (visité le 12/06/2022).

5. D. N. Smith et C. W. Blackwell, « Four URLs, limitless apps : Separation of concerns in the Homer Multitext architecture »...

6. B. Almas, T. Clérice, H. Cayless, *et al.*, *Distributed Text Services (DTS) : a Community-built API to Publish and Consume Text Collections as Linked Data...*

7. *Expression régulière*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Expression\\_r%C3%A9guli%C3%A8re](https://fr.wikipedia.org/wiki/Expression_r%C3%A9guli%C3%A8re) (visité le 10/07/2022).

8. M. Boeckhout, G. A. Zielhuis et A. L. Bredenoord, « The FAIR guiding principles for data stewardship... ».

**Fichier de log** Un fichier de log est un fichier qui est créé pendant l'exécution d'un algorithme pour sauvegarder des informations qui le concernent. Il s'agit en général d'informations annexes sur l'exécution du script : messages d'erreur, données produites en cours d'utilisation.... 71, 72, 76–78, 87, 88, 157

**HTTP** Le *Hypertext Transfer Protocol* (HTTP, ou P« protocole de transfert hypertexte »), est un protocole qui permet la communication entre différentes machines dans une architecture client-serveur, comme celle du Web. C'est donc le protocole qui est à la base d'Internet, puisqu'il établit la manière dont des machines peuvent interagir au sein d'un réseau. Il définit un ensemble de méthodes pour interagir avec le serveur – demander (GET) des données, en envoyer (POST)... Il définit également la manière dont le serveur répond au client et un ensemble d'erreurs<sup>9</sup>, chacune disposant d'un code pour l'identifier (l'erreur 404 indiquant que la ressource demandée n'existe pas, ou l'erreur 418 signifiant « Je suis une théière », un canular qui pousse l'humour des concepteurs du Web).. 67

**Named Entity Linking (NEL)** Le *Named Entity Linking*, traduit en français par « résolution d'entités nommées » ou « liage d'entités nommées » consiste à lier des entités – c'est-à-dire des objets uniques, comme des lieux, des personnes ou des organisations – à des sources de connaissance externes, comme des bases de connaissance en ligne.. 22

**OAI-PMH** L'OAI-PMH (*Open Archives Initiative Protocol for Metadata Harvesting*) est un protocole notamment utilisé par les bibliothèques et les archives ; il permet de partager des métadonnées sur une collection de documents. Les données sont partagées en XML et s'appuient fortement sur l'usage du *Dublin Core*, un standard simple pour la structuration de métadonnées. Contrairement au CTS et au DTS, l'API OAI-PMH ne permet pas de partager un document, mais seulement ses métadonnées ; comme ces deux autres standards cependant, le standard OAI-PMH ne définit pas une implémentation<sup>10</sup>.. 107

**Ontologie** Une ontologie correspond à la définition d'un ensemble de catégories, de propriétés et de relations qui unissent des données et des concepts ; cet ensemble est complété par une modélisation (souvent sous forme graphique), qui indique la relation entre les différents termes de l'ontologie<sup>11</sup>. Ceux-ci sont souvent liés de façon hiérarchique (plusieurs termes spécifiques pouvant être dérivés d'un terme générique). Dans le Web sémantique, chaque terme a lui-même un identifiant unique,

---

9. R. Fielding, *Architectural Styles and the Design of Network-based Software Architectures...*, p. 116-127.

10. C. Prime-Claverie et A. Mahé, « Le défi de l'interopérabilité entre plates-formes pour la construction de savoirs augmentés en sciences humaines et sociales »...

11. *Ontology (information science)*...

ce qui garantit une implémentation uniforme pour l'ontologie et réduit le risque de présence de doublons.. 20

**Python** Python est un langage de programmation impérative (c'est-à-dire, basé sur la production et la transformation de données par une suite d'instructions) créé en 1991. Il est particulièrement utilisé dans le domaine du traitement de données et des humanités numériques. C'est le langage le plus utilisé dans le projet *MSS / Katabase*. 44, 118

**REST** L'architecture REST, définie par Roy Fielding (2000) <sup>12</sup> un standard pour le design et le fonctionnement interne d'API. Contrairement à d'autres standards (DTS, CTS, OAI-PMH), elle ne définit pas la manière de cibler des données grâce à des identifiants; cette architecture définit plutôt comment l'interaction client-serveur doit être implémentée, à l'aide d'une série de principes de design. C'est le standard le plus répandu pour la conception d'API, bien qu'il soit souvent mal implémenté, ce qui a amené au développement d'API *REST-like*.. 67

**Score BLEU** Un score BLEU (bilingual evaluation understudy) est un score servant à évaluer la qualité de la traduction d'un texte par une machine. Il est obtenu en comparant le texte traduit par une machine avec un plusieurs textes de référence (c'est-à-dire, des traductions faites par des humains du même texte).. 128

**Score F1** Le score F1, ou *F-score*, est la moyenne harmonique de la précision (vrais positifs par rapport au total de résultats obtenus) et du rappel (nombre de résultats positifs par rapport au total de résultats positifs). <sup>13</sup>. Le score F1 a l'avantage de prendre en compte les vrais et les faux positifs. Ce score, dont la valeur est contenue entre 0 et 1, permet de mesurer l'exactitude d'un algorithme d'apprentissage machine, ou d'un moteur de recherche. Soit *F* le score F1, *P* la précision et *R* le rappel, sa formule est la suivante :

$$P = \frac{\text{items pertinents obtenus}}{\text{total des items obtenus}}$$

$$R = \frac{\text{items pertinents obtenus}}{\text{ensemble des items pertinents}}$$

$$F = 2 \times \frac{P \times R}{P + R}$$

. 25, 78, 79, 157

---

12. R. Fielding, *Architectural Styles and the Design of Network-based Software Architectures...*

13. *Précision et rappel*, Wikipedia. L'encyclopédie libre, 2022, URL : [https://fr.wikipedia.org/wiki/Pr%C3%A9cision\\_et\\_rappel](https://fr.wikipedia.org/wiki/Pr%C3%A9cision_et_rappel) (visité le 13/06/2022); *F-score*, Wikipedia. L'encyclopédie libre, 2022, URL : <https://en.wikipedia.org/wiki/F-score> (visité le 12/06/2022).



**SPARQL** SPARQL est un langage de requête qui permet d’interagir avec une base de données au format RDF. « SPARQL exprime des requêtes sur des sources de données diverses [...]. SPARQL rend possible la requête de données en graphes [...], avec des données conjointes et disjointes. SPARQL permet également l’agrégation de données, les sous-requêtes, la négation, la création de données à l’aide d’expression, le test de données et la contrainte des requêtes. Les résultats de requêtes sparql peuvent être des jeux de résultats ou des graphes RDF. »<sup>14</sup>. 13, 20, 21, 23–27, 65, 81–94, 96, 98, 122, 137, 145, 151, 185, 189, 196

**Table de conversion** Une table de conversion est, tout simplement, un dictionnaire qui contient en clés un certain nombre d’informations telles qu’elles figurent dans le texte et, en valeurs, une version normalisée de cette information. Cela permet de détecter des motifs à extraire autant de normaliser les informations.. 25, 35, 39

**URI** Une URI (*Uniform Resource Identifier*) est une chaîne de caractères qui identifie une ressource sur un réseau. Cette URI est permanente, et permettra d’identifier une ressource même si celle-ci est supprimée. L’URL est un type particulier d’URI qui, en plus de décrire une ressource, la représente aussi : un URL de page *Wikipedia* permet d’identifier une page de façon unique, comme toute URI (aucune autre page sur internet n’ayant le même URL) ; mais en plus, elle permet d’accéder à la page elle-même.. 21

**URL** Une URL (*Uniform Resource Locator*) est un identifiant unique sous la forme d’une chaîne de caractères pointant vers une ressource sur le Web. De la même manière qu’une cote de bibliothèque pointe vers un livre précis, un URL pointe vers une page ou une donnée sur le Web. L’URL permet d’agir sur une ressource (en envoyant des données...) et peut être associée à des données ainsi qu’à des représentations (une page Web, par exemple).. 69

---

14. « SPARQL can be used to express queries across diverse data sources [...]. SPARQL contains capabilities for querying [...] graph patterns along with their conjunctions and disjunctions. SPARQL also supports aggregation, subqueries, negation, creating values by expressions, extensible value testing, and constraining queries by source RDF graph. The results of SPARQL queries can be result sets or RDF graphs. » Steve Harris, Andy Seaborne et Eric Prud’hommeaux, *SPARQL 1.1 Query Language. W3C Recommendation 21 March 2013*, W3C, 2013, URL : <http://www.w3.org/TR/2013/REC-sparql11-query-20130321/> (visité le 10/06/2022) (traduction de l’auteur).



# Acronymes

**API** *Application Programming Interface*. 20, 22, 23, 25, 30, 32, 33, 36, 37, 40, 45, 53, 56, 57, 59, 61, 65–72, 74–76, 78, 99, 105–107, 109–123, 125, 145, 150, 151, 183, *Glossaire* : API

**CTS** *Canonical Text Services*. 107, 113–117, *Glossaire* : CTS

**DTS** *Distributed Text Services*. 107, 113, 116, 117, *Glossaire* : DTS

**FAIR** *Findable Accessible Interoperable Reusable*. 107–109, 112, 119, *Glossaire* : FAIR

**HTTP** *Hypertext Transfer Protocol*. 67–69, 87, 105, 110, 111, 121, 122, *Glossaire* : HTTP

**NEL** *Named Entity Linking (NEL)*. 22, 23, 78, 79, 99, 147, 159, *Glossaire* : Named Entity Linking (NEL)

**OAI-PMH** *Open Archives Initiative Protocol for Metadata Harvesting*. 107, 113, 114, 117, *Glossaire* : OAI-PMH

**REST** *Representational state transfer*. 67, 68, 107, 110, 112, 114, 119, 120, 122, *Glossaire* : REST

**URI** *Uniform Resource Identifier*. 21, 90, 91, 110, 115, 116, *Glossaire* : URI

**URL** *Uniform Resource Locator*. 69–72, 96, 98, 105–107, 110, 115, 117, *Glossaire* : URL



# Table des figures

3.1	Deux exemples de lettres . . . . .	18
3.2	Présentation générale de l'algorithme d'enrichissement de données à l'aide de <i>Wikidata</i> . . . . .	27
4.1	Processus d'extraction d'informations du <code>tei:name</code> et <code>tei:item</code> . . . . .	41
4.2	Les différentes parties du <code>tei:name</code> . . . . .	47
4.3	Différentes structures possibles pour un <code>tei:name</code> . . . . .	49
4.4	Extraction d'informations d'un <code>tei:name</code> contenant des parenthèses . . . .	52
4.5	Représentation graphique du processus d'extraction et de reconstitution d'un nom à partir de son abréviation . . . . .	58
5.1	L'interaction client-serveur au travers d'une API . . . . .	68
5.2	Résultats retournés par <i>Wikidata</i> pour la recherche « Virginia Woolf » . . .	70
5.3	Algorithme lançant des recherches en plein texte sur <i>Wikidata</i> . . . . .	74
6.1	Exemple de relation sujet – prédicat – objet . . . . .	82
6.2	L'algorithme de constitution d'une base de connaissances issues de <i>Wiki-</i> <i>data</i> à l'aide de SPARQL . . . . .	88
1	. . . . .	162
2	Répartition des différents types de <code>tei:name</code> . . . . .	163
3	De la TEI à l'API : document en entrée, informations extraites et chaînes de caractères recherchées . . . . .	183
4	L'architecture des bases de données de la <i>Wikimedia Foundation</i> (schéma réalisé par Timo Tijhof en 2020, disponible dans le domaine public avec la licence Creative Commons CC0.1.0 Universal Public Domain License) . . .	200



# Liste des tableaux

1	Tests menés avec un jeu de 200 entrées sur des paramètres isolés . . . . .	156
2	Tests menés avec un jeu de 200 entrées sur l’algorithme final . . . . .	157
3	Résultats du NEL de lieux dans des textes littéraires du XIXe s. par Sou- dani et al. (2018) . . . . .	159





# Table des matières

Résumé	i
Introduction	1
<b>I Du document numérisé au XML-TEI : nature du corpus, structure des documents et méthode de production des données</b>	<b>3</b>
1 Le marché des manuscrits autographes au prisme des catalogues de vente	5
1.1 Pourquoi étudier le marché des manuscrits autographes? . . . . .	5
1.2 La structure du corpus : périodisation, producteurs des documents et classification . . . . .	5
1.2.1 Le corpus de catalogues de vente de manuscrits . . . . .	6
1.2.2 Structure des catalogues . . . . .	6
2 Production des données : de l'OCR à la TEI	7
2.1 Extraire le texte des imprimés . . . . .	7
2.1.1 Comprendre la structure du document pour préparer l'édition numérique . . . . .	7
2.2 L'encodage des manuscrits en XML-TEI . . . . .	8
2.2.1 Encoder les catalogues en TEI . . . . .	8
2.2.2 L'encodage en TEI : un processus sélectif qui réduit les significations du texte . . . . .	8
<b>II Vers une étude des facteurs déterminant le prix des documents : alignement des entrées de catalogue et enrichissement sémantique avec <i>Wikidata</i></b>	<b>11</b>
3 Questions introductives : pourquoi et comment s'aligner avec <i>Wikidata</i> ?	15
3.1 Questions théoriques . . . . .	15

3.1.1	Pourquoi s'aligner avec des entités externes ? . . . . .	15
3.1.2	Pourquoi utiliser la base de connaissances <i>Wikidata</i> ? . . . . .	19
3.1.3	Quelle relation avec la résolution d'entités nommées d'entités nom- mées ? . . . . .	22
3.2	Questions de méthode : comment faire le lien entre un corpus historique et une base de connaissances en ligne ? . . . . .	24
3.2.1	Présentation générale de l'algorithme . . . . .	24
3.2.2	Comment traduire des descriptions textuelles datant du XIX <sup>ème</sup> s. en chaînes de caractères pouvant retourner un résultat sur <i>Wikidata</i> ?	28
3.2.3	Comment négocier avec le moteur de recherche de <i>Wikidata</i> ? . . .	32
3.3	Une approche prédictive . . . . .	34
<b>4</b>	<b>Un algorithme de détection de motifs pour préparer et structurer les données</b>	<b>37</b>
4.1	Présentation générale . . . . .	37
4.1.1	Les formats d'entrée et de sortie . . . . .	37
4.1.2	Présentation de l'algorithme d'extraction d'informations . . . . .	39
4.2	Identifier le type de nom . . . . .	42
4.3	Le traitement des noms de personnes . . . . .	45
4.3.1	Trouver des solutions adaptées à différents types de noms . . . . .	46
4.3.2	Identifier et extraire les informations nominatives . . . . .	48
4.3.3	Reconstruire un prénom complet à partir de son abréviation . . . . .	53
4.4	Extraire des informations biographiques du <code>tei:trait</code> . . . . .	59
4.4.1	Identifier les dates de vie et de mort d'une personne . . . . .	59
4.4.2	Identifier l'occupation d'une personne . . . . .	61
<b>5</b>	<b>Résoudre les entités par l'alignement avec des identifiants <i>Wikidata</i></b>	<b>65</b>
5.1	Quantifier l'incertitude : quels sont les paramètres qui permettent d'obtenir le bon résultat ? . . . . .	65
5.2	Automatiser l'extraction de données depuis des sources externes : les API .	67
5.2.1	Qu'est-ce qu'une API ? . . . . .	67
5.2.2	Consommer une API : le cas de l'application <i>MediaWiki</i> . . . . .	69
5.3	Présentation générale . . . . .	71
5.4	Gérer la montée en charge : optimisation et réduction du temps d'exécution	75
5.5	Évaluer la résolution d'entités : performance, qualité des données extraites de <i>Wikidata</i> et comparaison avec d'autres projets . . . . .	77
5.6	En conclusion : retour sur l'extraction d'informations des catalogues et sur l'algorithme d'alignement avec <i>Wikidata</i> . . . . .	79

<b>6</b>	<b>Après l'alignement, l'enrichissement : utiliser SPARQL pour produire des données structurées</b>	<b>81</b>
6.1	Comprendre les particularités des modèles sémantiques de données . . . . .	81
6.2	Quelles données rechercher via SPARQL ? . . . . .	83
6.3	Présentation générale . . . . .	86
6.4	Développer un comportement uniforme pour produire des données exploitables à partir un corpus hétérogène . . . . .	89
6.4.1	Les spécificités techniques de SPARQL : adapter les résultats à ses besoins . . . . .	89
6.4.2	La base de connaissances à construire . . . . .	91
6.4.3	Des réponses SPARQL à la base de connaissance . . . . .	93
<b>7</b>	<b>Un corpus augmenté : mise en relation du corpus aux données de <i>Wikidata</i> et possibilités ouvertes par cet alignement</b>	<b>95</b>
7.1	Lier la TEI aux données nouvellement produites . . . . .	95
7.2	Que faire des entités ? . . . . .	96
7.3	En conclusion . . . . .	99
<b>III</b>	<b>Rendre la recherche réutilisable et interopérable : <i>KatAPI</i>, une API pour échanger des données structurées</b>	<b>103</b>
<b>8</b>	<b>Standards de design et statut des API dans pour les humanités numériques centrées sur le texte</b>	<b>107</b>
8.1	Que faire du FAIR ? Le partage des données de la recherche . . . . .	108
8.2	Le REST : un modèle pour le design d'API . . . . .	109
8.3	OAI-PMH, CTS et DTS : quels standards pour le partage du texte en humanités numériques ? . . . . .	112
8.3.1	OAI-PMH : un premier standard à succès . . . . .	113
8.3.2	CTS, DTS : des méthodes de standardisation pour l'échange de texte	114
8.4	Pour qui sont ces API ? Qu'est-ce qui est demandé d'un tel outil ? . . . . .	117
<b>9</b>	<b>Définir un périmètre : que partager, et comment partager ?</b>	<b>121</b>
9.1	Grands principes pour l'architecture de l'API . . . . .	121
9.2	Quelles données partager ? . . . . .	122
9.3	Codifier l'accès aux données : présentation des paramètres de recherche possibles . . . . .	123
9.4	Comment partager les données ? Principes suivis pour le partage d'informations, formats et structure des réponses de l'API . . . . .	123

<b>10 Implémentation et fonctionnement technique de <i>KatAPI</i></b>	<b>125</b>
10.1 Présentation générale . . . . .	125
10.2 Gestion des erreurs . . . . .	125
10.3 Garantir le bon fonctionnement de l'application . . . . .	125
<b>Conclusion</b>	<b>127</b>
<b>Bibliographie</b>	<b>129</b>
Projet <i>MSS</i> / <i>Katabase</i> . . . . .	129
Édition numérique, traitement automatisé et analyse de texte . . . . .	130
Technologies du web : généralités . . . . .	132
Technologies du web : Web sémantique . . . . .	133
Technologies du web : API . . . . .	133
Visualisation et design d'interfaces . . . . .	135
Marché de l'art, économétrie et statistiques . . . . .	136
<b>Glossaire</b>	<b>137</b>
<b>Acronymes</b>	<b>143</b>
<b>Table des figures</b>	<b>145</b>
<b>Liste des tableaux</b>	<b>147</b>
<b>Table des matières</b>	<b>149</b>
<b>Lien vers tout le code source de la chaîne de traitement <i>Katabase</i></b>	<b>153</b>
<b>Résultat des tests de l'algorithme d'extraction d'informations de <i>Wikidata</i></b>	<b>155</b>
<b>Graphiques</b>	<b>161</b>
<b>Code source et données encodées</b>	<b>165</b>
<b>Images</b>	<b>199</b>

# Lien vers tout le code source de la chaîne de traitement *Katabase*

L'intégralité du code produit pour le projet *Katabase*, depuis la création des documents TEI jusqu'à l'application Web, sont disponibles sur la plateforme d'hébergement de code GitHub.

- Plateforme du projet
- Première étape : normalisation des documents TEI
- Deuxième étape : extraction de données et normalisation
- Troisième étape : enrichissement de données à l'aide de *Wikidata*
- Quatrième étape : création de jeux de données JSON pour le site web
- Application web *Katabase*

L'intégralité du code source est disponible sous licence libre (GNU GPL v3.0, MIT ou Creative Commons). Le site web peut également être consulté en ligne à cette adresse.



# Résultat des tests de l'algorithme d'extraction d'informations de *Wikidata*

Les résultats des tests présentés dans les tables ci-dessous (1, 2) ont été menés sur un jeu de 200 couples `tei:name` et `tei:trait`. Ils ont été choisis de façon à être représentatifs du jeu de données complet, avec une variété d'entrées semblables (personnes nobles et non nobles, entrées qui ne sont pas consacrées à des personnes...). Ce jeu de test présente également une structure et un degré de détail semblable au jeu de données complet. Il y a notamment exactement la même proportion d'entrées présentant un `tei:trait` que dans le jeu de test et le jeu final.

Les résultats en pourcentage sont toujours exprimés en **pourcentages du nombre total d'entrées requêtées** (soit 200).

Type de requête	Noms	Noms et de nom famille noble	Noms et de titre noblesse	Noms et date de naissance et mort	Noms et de occupation
Avec reconstitution des prénoms	48%	45.9%	50%	52.6%	53.1%
Sans reconstitution des prénoms	42%	40.5%	50%	52.6%	53.1%

TABLE 1 – Tests menés avec un jeu de 200 entrées sur des paramètres isolés

Le tableau ci-dessus (1) présente les résultats de tests menés pour isoler l'impact de chaque paramètre dans l'obtention du bon résultat. Les mêmes requêtes ont été lancées sur un jeu de test de 200 entrées de catalogue. Toutes les requêtes contiennent le prénom et le nom de famille complets (« Noms » dans la table). La première requête ne contient que ce paramètre, tandis que toutes les autres requêtes sont lancées avec un paramètre en plus : titre de noblesse, nom de famille noble, dates et occupation. Le fait que ces paramètres soient utilisés ne garantit pas qu'ils soient disponibles dans le jeu de données utilisé en entrée. Ce test permet d'isoler de mieux construire l'algorithme final de recherches en plein texte sur le moteur de recherche de *Wikidata*, puisqu'il permet de voir quels sont les paramètres les plus fiables.



<b>Bons identifiants récupérés</b>	65%
<b>Score F1</b>	0.674
<b>Précision</b>	0.677
<b>Rappel</b>	0.670
<b>Résultats certains</b>	36%
<b>Faux positifs parmi les éléments certains</b>	8.5%
<b>Total d'entrées requêtées</b>	200
<b>Identifiants récupérés</b>	192
<b>Silence</b>	8
<b>Temps d'exécution avec l'option <i>fetch</i></b>	88.3 secondes
<b>Temps d'exécution sans l'option <i>fetch</i></b>	92.49 secondes

TABLE 2 – Tests menés avec un jeu de 200 entrées sur l'algorithme final

Le tableau ci-dessus présente les résultats d'un test évaluant les performances de l'algorithme final. Voici une explication de ces résultats :

- *Bons identifiants récupérés* : la proportion, exprimée en pourcentage, du nombre d'identifiants récupérés par l'algorithme qui correspondent aux identifiants récupérés manuellement
- *Score F1* : le score F1, soit la moyenne pondérée de la précision et du rappel. Contrairement à la mesure ci-dessus, le score F1 prend en compte le silence et les faux négatifs, et offre donc une interprétation plus complète de la qualité du script.
- *Précision* : la précision est une composante du score F1 qui mesure la proportion d'identifiants pertinents obtenus parmi tous les identifiants obtenus.
- *Rappel* : le rappel est une composante du score F1 qui mesure la proportion d'identifiants pertinents obtenus parmi l'ensemble des identifiants pertinents qui existent.
- *Résultats certains* : pour accélérer le processus de relecture, un score de certitude a été attribué à certains identifiants obtenus grâce à l'algorithme. Ce score dépend du nombre de paramètres utilisés dans la recherche qui a permis de récupérer l'identifiant. 36% des résultats obtenus sont considérés comme certains.
- *Faux positifs parmi les éléments certains* : le score de certitude présenté ci-dessus admet un taux d'erreur qui est ici quantifié : 8.5% des résultats obtenus sont erronés alors qu'ils ont été considérés comme certains (8,5% du jeu de données total, donc).
- *Total d'entrées requêtées* : le volume du jeu de données de test, soit 200 entrées.
- *Identifiants récupérés* : le nombre d'identifiants obtenus, qu'ils soient corrects ou non.
- *Silence* : le nombre d'entrées pour lesquelles aucun identifiant n'a été obtenu.
- *Temps d'exécution avec l'option *fetch** : le temps pris, en secondes, pour lancer l'intégralité de l'algorithme de récupération d'identifiants *Wikidata* sur les 200 entrées en stockant dans des fichiers de log les requêtes lancées et les identifiants requêtés.

- *Temps d'exécution sans l'option **fetch*** : le temps pris, en secondes, pour lancer l'intégralité de l'algorithme de récupération d'identifiants sur les 200 entrées sans utiliser de fichiers de log.

	Base de connaissances	Score F1 – recherche des candidates	Score F1 – Lien avec le candidat	Overall linking accuracy
Score				
DBPedia		0.899	0.536	0.834
DataBNF		0.689	0.726	0.7
Wikidata		0.869	0.611	0.85

TABLE 3 – Résultats du NEL de lieux dans des textes littéraires du XIXe s. par Soudani et al. (2018)

Scores obtenus par Aicha Soudani et al. pour le liage d’entités nommées pour des lieux dans des textes littéraires français du XIX<sup>ème</sup> s., présentés à la conférence *Hu-manNS’2018*<sup>15</sup>. Les scores F1 ne figurent pas dans l’article original, qui ne présente que la précision et le rappel. Je les ai donc calculés moi même. La méthode de calcul de l’*overall linking accuracy* (« exactitude générale du liage ») n’est pas précisée dans l’article ; il est seulement indiqué que « Cette mesure essaie d’évaluer l’efficacité globale du système, et non par phase, et exprime la fiabilité de REDEN pour la tâche de résolution des entités nommées. »<sup>16</sup>.

La méthode utilisée dans cet article est la suivante : les auteur.ice.s s’appuient fortement sur l’apprentissage machine, en identifiant d’abord des entités dans le texte à l’aide de SEM avant de lier les entités à différentes bases de connaissance en utilisant REDEN.

---

15. A. Soudani, Y. Meherzi, A. Bouhafs, *et al.*, « Adaptation et évaluation de systèmes de reconnaissance et de résolution des entités nommées pour le cas de textes littéraires français du 19ème siècle »..., p. 4.

16. *Ibid.*



# Graphiques

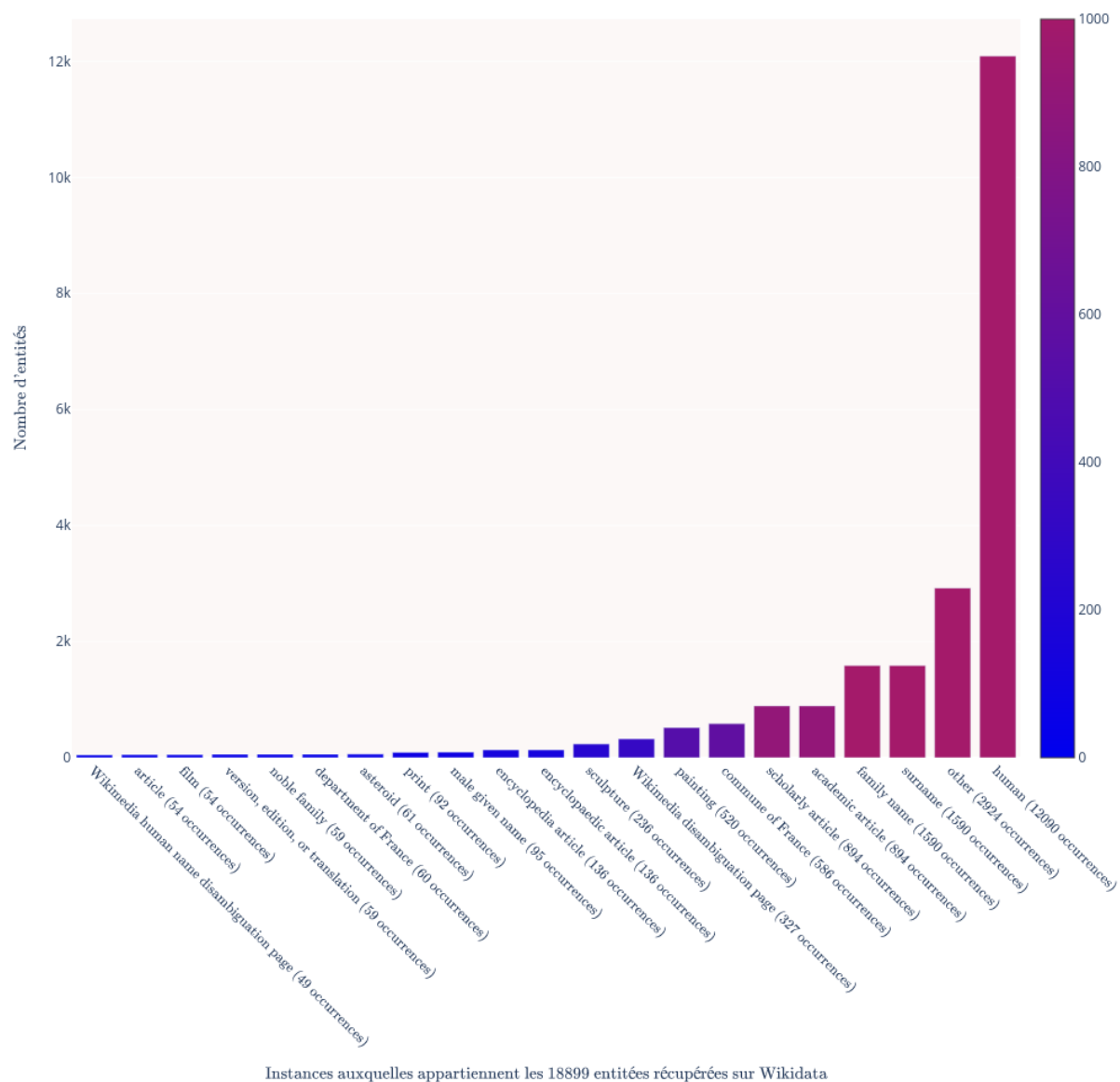
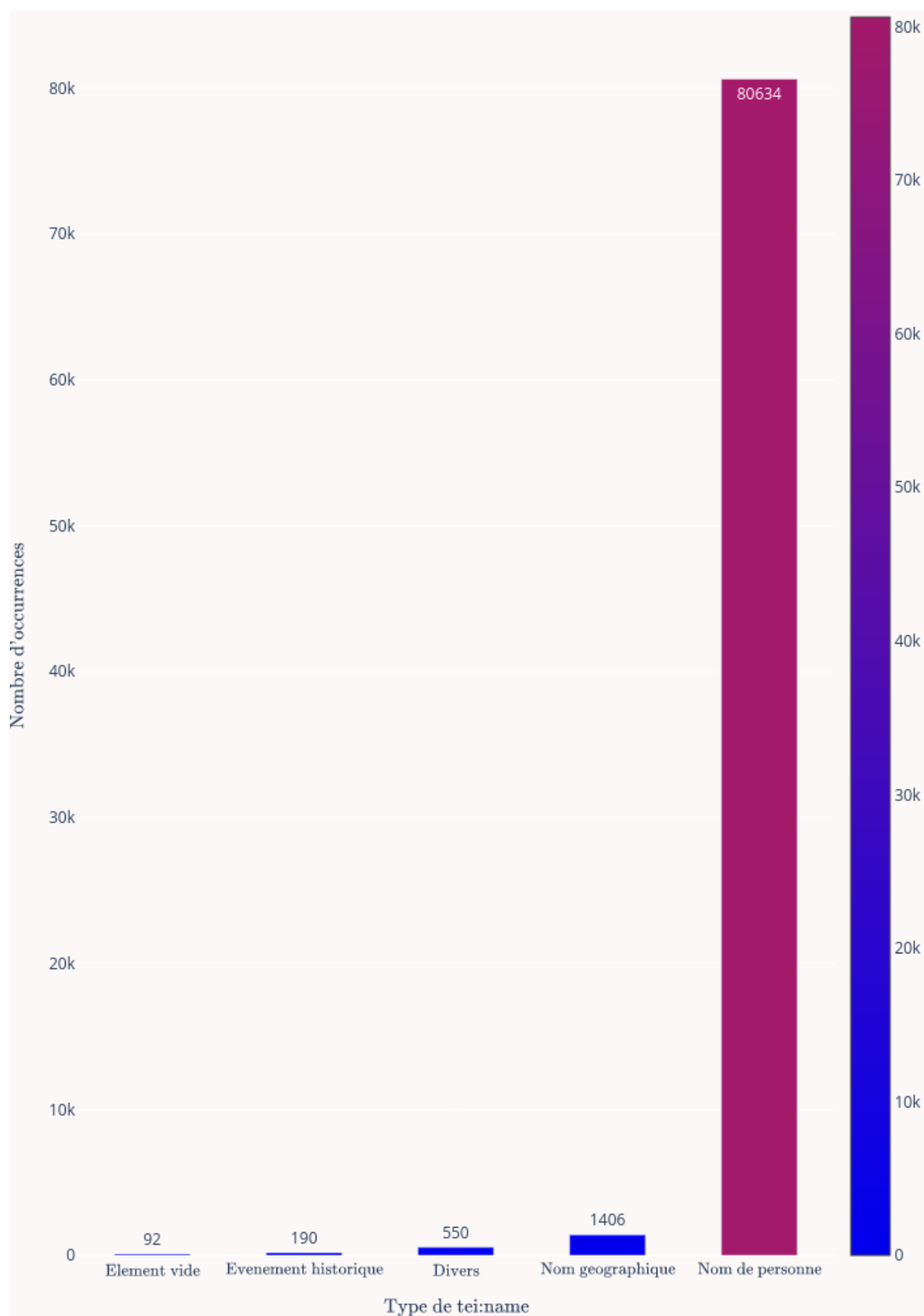


FIGURE 1 –  
Occurrences des différentes catégories auxquelles appartiennent les entités *Wikidata* liées  
avec les entrées de catalogues

FIGURE 2 – Répartition des différents types de `tei:name`





## Code source et données encodées

```

1  functions = {
2      "général": "general",
3      "maréchal": "marshal",
4      "lieutenant": "military",
5      "officier": "military",
6      "colonel": "military",
7      "lieutenant-colonel": "military",
8      "commandant": "military",
9      "capitaine": "military",  # "less important" military positions
10     "roi": "king",
11     "empereur": "emperor",
12     "président": "president",
13     "homme politique": "politician",
14     "président de l'assemblée": "politician",
15     "orateur": "politician",
16     "député": "politician",
17     "secrétaire d'état": "politician",
18     "sénateur": "politician",
19     "écrivain": "writer",
20     "auteur": "writer",
21     "romancier": "writer",
22     "acteur": "actor",
23     "actrice": "actress",
24     "cantatrice": "singer",
25     "chanteur": "singer",
26     "chanteuse": "singer",
27     "peintre": "painter",
28     "sculpteur": "sculptor",
29     "statutaire": "sculptor",
30     "compositeur": "composer",
31     "musicien": "musician",
32     "musicienne": "musician",
33     "tragédien": "actor",
34     "chansonnier": "chansonnier",
35     "architecte": "architect",
36     "journaliste": "journalist",
37     "inventeur": "inventor",
38     "chimiste": "chemist",
39     "connétable": "constable",
40     "archevêque": "archbishop",
41     "évêque": "bishop",
42     "docteur": "physicist",
43     "médecin": "physicist"
44 }

```

Code source 1 – Table de conversion associant un métier à son équivalent normalisé

```
1 dpts = [  
2     "ain",  
3     "aisne",  
4     "allier",  
5     "basses-alpes",  
6     "hautes-alpes",  
7     "alpes-maritimes",  
8     "annepins",  
9     "provence",  
10    "ardèche",  
11    "ardennes",  
12    "arriège",  
13    "arno",  
14    "aube",  
15    "aude",  
16    "aveyron",  
17    "bouches-de-l'elbe",  
18    "bouches-de-l'escaut",  
19    "bouches-de-l'yssel",  
20    "bouches-de-la-meuse",  
21    "bouches-du-rhin",  
22    "bouches-du-rhône",  
23    "bouches-du-weser",  
24    "calvados",  
25    "cantal",  
26    "charente",  
27    "charente-inférieure",  
28    "cher",  
29    "corrèze",  
30    "corse",  
31    "côte-d'or",  
32    "côtes-du-nord",  
33    "creuse",  
34    "deux-nèthes",  
35    "deux-sèvres",  
36    "doire",  
37    "dordogne",
```

```
38     "doubs",
39     "drôme",
40     "dyle",
41     "ems-occidental",
42     "ems-oriental",
43     "ems-supérieur",
44     "escaut",
45     "eure",
46     "eure-et-loir",
47     "finistère",
48     "forêts",
49     "gard",
50     "haute-garonne",
51     "gers",
52     "gironde",
53     "hérault",
54     "ille-et-villaine",
55     "indre",
56     "indre-et-loire",
57     "isère",
58     "jemappes",
59     "jura",
60     "landes",
61     "léman",
62     "loire",
63     "loir-et-cher",
64     "haute-loire",
65     "loire-inférieure",
66     "loiret",
67     "lot",
68     "lot-et-garonne",
69     "lozère",
70     "lys",
71     "maine-et-loire",
72     "manche",
73     "marengo",
74     "marne",
75     "haute-marne",
```

76 "méditerranée",  
77 "mayenne",  
78 "meurthe",  
79 "meuse",  
80 "meuse-inférieure",  
81 "mont-blanc",  
82 "mont-tonnerre",  
83 "montenotte",  
84 "morbihan",  
85 "meuse",  
86 "moselle",  
87 "nièvre",  
88 "nord",  
89 "oise",  
90 "ombrone",  
91 "orne",  
92 "ourte",  
93 "paris",  
94 "pas-de-calais",  
95 "pô",  
96 "puy-de-dôme",  
97 "hautes-pyrénées",  
98 "basses-pyrénées",  
99 "pyrénées-orientales",  
100 "haut-rhin",  
101 "bas-rhin",  
102 "rhin-et-moselle",  
103 "rhône",  
104 "rhône-et-loire",  
105 "roer",  
106 "rome",  
107 "haute-saône",  
108 "saône-et-loire",  
109 "sambre-et-meuse",  
110 "sarre",  
111 "sarthe",  
112 "seine",  
113 "seine-et-marne",

```
114     "seine-et-oise",
115     "seine-inférieure",
116     "sézia",
117     "simplon",
118     "deux-sèvres",
119     "somme",
120     "stura",
121     "tarn",
122     "tarn-et-garonne",
123     "taro",
124     "trasimène",
125     "var",
126     "vaucluse",
127     "vendée",
128     "vienne",
129     "haute-vienne",
130     "vosges",
131     "yonne",
132     "yssel-supérieur",
133     "zuyderzée"
134 ]
```

Code source 2 – Liste de départements du XIXe s. pour détecter des informations géographiques

```

1 countries = {
2     "états-unis d'amérique": "united states of america",
3     "etats-unis d'amérique": "united states of america",
4     "états unis d'amérique": "united states of america",
5     "etats unis d'amerique": "united states of america",
6     "états-unis": "united states of america",
7     "etats-unis": "united states of america",
8     "etats unis": "united states of america",
9     "états unis": "united states of america",
10    "italie": "italy",
11    "grèce": "greece",
12    "canada": "canada",
13    "chine": "china",
14    "haïti": "haiti",
15    "tobago": "tobago",
16    "brésil": "brasil",
17    "burkina-faso": "burkina-faso",
18    "cameroun": "cameroun",
19    "tchad": "tchad",
20    "congo": "congo",
21    "gabon": "gabon",
22    "guinée": "guinea",
23    "côte d'ivoire": "ivory coast",
24    "mali": "mali",
25    "mauritanie": "mauritania",
26    "niger": "niger",
27    "sénégal": "senegal",
28    "madagascar": "madagascar",
29    "seychelles": "seychelles",
30    "tanzanie": "tanzania",
31    "zanzibar": "zanzibar",
32    "liban": "lebanon",
33    "syrie": "syria",
34    "inde": "india",
35    "laos": "laos",
36    "viet-nâm": "vietnam"
37 }

```

Code source 3 – Table de conversion pour les pays

```
1 colonies = [  
2     "québec",  
3     "ontario",  
4     "saint-pierre-et-miquelon",  
5     "mississippi",  
6     "missouri",  
7     "louisiane",  
8     "anguilla",  
9     "antigua",  
10    "dominique",  
11    "saint-domingue",  
12    "guadeloupe",  
13    "monsterrat",  
14    "saint-martin",  
15    "saint-barthélémy",  
16    "sainte-lucy",  
17    "saint-vincent-et-les-grenadines",  
18    "saint-eustache",  
19    "saint-christophe",  
20    "martinique"  
21    "guyane française",  
22    "guyane",  
23    "maroc", # unfortunately the morocco referred to in XIXth century  
24    ↪ france is a french protectorate  
25    "algérie", # same  
26    "algérie française", # same  
27    "tunisie", # same  
28    "fezzan",  
29    "dahomey",  
30    "haute-volta",  
31    "oubangui-chari",  
32    "congo français",  
33    "moyen-congo",  
34    "guinée française",  
35    "soudan français",  
36    "gorée",  
    "tigi",
```



```

37     "djibouti",
38     "cheikh saïd",
39     "comores",
40     "fort-dauphin",
41     "îles maurice",
42     "mayotte",
43     "la réunion",
44     "îles éparses",
45     "île amsterdam",
46     "île saint-paul",
47     "archipel crozet",
48     "îles kerguelen",
49     "castellorizo",
50     "grand-liban",
51     "sandjak d'alexandrette",
52     "indes françaises",
53     "pondichéry",
54     "karikal",
55     "yanaon",
56     "mahé",
57     "chanderngor",
58     "tonkin",
59     "annam",
60     "cochinchine",
61     "guangzhou wan",
62     "shanghai",
63     "guangzhou",
64     "tianjin",
65     "hankou",
66     "clipperton",
67     "nouvelle-calédonie",
68     "polynésie française",
69     "vanuatu",
70     "nouvelles-hébrides",
71     "wallis et futuna"
72 ]

```

Code source 4 – Liste d'anciennes colonies françaises utilisées pour la détection de motifs

```
1 provinces = [  
2     "armagnac",  
3     "île-de-france",  
4     "berry",  
5     "orléanais",  
6     "normandie",  
7     "languedoc",  
8     "lyonnais",  
9     "dauphiné",  
10    "champagne",  
11    "aunis",  
12    "saintonge",  
13    "poitou",  
14    "guyenne et gascogne",  
15    "bourgogne",  
16    "picardie",  
17    "anjou",  
18    "provence",  
19    "angoumois",  
20    "bourbonnais",  
21    "marche",  
22    "bretagne",  
23    "maine",  
24    "touraine",  
25    "limousin",  
26    "comté de foix",  
27    "auvergne",  
28    "béarn",  
29    "alsace",  
30    "artois",  
31    "roussillon",  
32    "flandre française et hainaut français",  
33    "franche-comté",  
34    "lorraine et trois-évêchés",  
35    "corse",  
36    "nivernais",  
37 ]
```

Code source 5 – Liste d'anciennes provinces françaises pour la détection de motifs

```
1 events = {  
2     "défense nationale": "government of national defense",  
3     "defense nationale": "government of national defense",  
4     "révolution française": "french revolution",  
5     "revolution francaise": "french revolution",  
6     "guerre de trente ans": "thirty years' war 1618 1648",  
7     "guerre de cent ans": "hundred years' war 1337 1453",  
8     "guerre de sept ans": "seven years war 1756 1763",  
9     "guerre": "war",  
10    "insurrection": "war",  
11    "siège de mayence": "siege of mainz",  
12    "siège": "siege",  
13    "commune": "commune",  
14    "défense": "battle",  
15    "révolution": "revolution"  
16 }
```

Code source 6 – Table de conversion pour les évènements historiques

```
1 status = {
2     "empereur": "",
3     "impératrice": "",
4     "général": "general",
5     "reine": "queen",
6     "roi": "king",
7     "princesse": "princess",
8     "prince": "prince",
9     "archiduchesse": "",
10    "archiduc": "",
11    "duchesse": "duchess",
12    "duc": "duke",
13    "famille": "family",
14    "seigneur": "",
15    "vicomtesse": "",
16    "victesse": "",
17    "vicomte": "",
18    "victe": "",
19    "comtesse palatine": "countess palatine",
20    "comtesse": "",
21    "ctesse": "",
22    "comte": "",
23    "cte": "",
24    "cardinal": "",
25    "pape": "pope",
26    "lord": "",
27    "chevalier": "",
28    "marquise": "",
29    "marquis": "",
30    "sire": "",
31    "baronnesse": "",
32    "baronne": "",
33    "baron": "",
34    "abbé": "",
35    "madame": "",
36    "mme": "",
37    "monsieur": "",
38    "mr": "",
39    "docteur": "",
40    "maréchale": "",
41    "maréchal": "",
42    "mademoiselle": "",
43    "melle": "",
44    "mlle": "",
45    "sir": ""
46 }
```

Code source 7 – Liste de colonies pour la détection de motifs

```

1 comp_names = {
2     "arm ch": "armand-charles",
3     "ch m": "charles-marie",
4     "ch l f": "charles-louis-françois",
5     "f m": "francois-marie",
6     "fr emm.": "françois-emmanuel",
7     "j ant": "jean-antoine",
8     "j f": "jean-francois",
9     "j m": "jean-marie",
10    "j j": "jean-jacques",
11    "j l": "jean-louis",
12    "j b": "jean-baptiste",
13    "j p": "jean-pierre",
14    "j pierre": "jean-pierre",
15    "l f": "louis-françois",
16    "m f": "marius-felix",
17    "franc rené": "francois-rené",
18    "m madeleine": "marie-madeleine",
19    "ph h": "philippe henri",
20    "p aug": "pierre auguste",
21    "p alex": "pierre alexandre",
22    "p j": "pierre-jean",
23    "j sylvain": "jean-sylvain",
24    "l ph": "louis-philippe",
25    "edm ch": "edmond-charles",
26    "ch marie": "charles-marie"
27 }

```

Code source 8 – Table de conversion permettant de remplacer un nom abrégé composé par sa version complète

```
1 names = {  
2     "ad": "adam",  
3     "alex": "alexandre",  
4     "alph": "alphonse",  
5     "ant": "antoine",  
6     "arm": "armand",  
7     "aug": "auguste",  
8     "ch": "charles",  
9     "cl": "claire",  
10    "dom": "dominique",  
11    "emm": "emmanuel",  
12    "ed": "edouard",  
13    "et": "etienne",  
14    "ét": "etienne",  
15    "ferd": "ferdinand",  
16    "fred": "frederic",  
17    "fr": "françois",  
18    "franc": "françois",  
19    "franç": "françois",  
20    "fréd": "frédéric",  
21    "g": "guillaume",  
22    "guill": "guillaume",  
23    "gab": "gabriel",  
24    "jh": "joseph",  
25    "jacq": "jacques",  
26    "jos": "joseph",  
27    "math": "matthieu",  
28    "nic": "nicolas",  
29    "ph": "philippe",  
30    "v": "victor",  
31    "vr": "victor",  
32 }
```

Code source 9 – Table de conversion permettant de remplacer un nom abrégé non composé par sa version complète

```

1 def rgx_abvcomp(nstr):
2     """
3     try to extract an abbreviated composed first name. if there is no
↳ match, return None
4     pattern
5     -----
6     the patterns in the example below are simplified to keep things
↳ readable
7     - two strings separated by a "-" or "\s"
8     - the first or second string can be a full name ([A-Z][a-z]+)
9     or an abbreviation ([A-Z][a-z]*\.)
10    - if the strings are separated by "\s", they must be finished by
↳ "\."
11    (to be sure that we don't capture full names, i.e: "J. Ch." can
↳ be captured,
12    but not "Jean Charles")
13    - complex names with 3 or more words must have "-" and at least
↳ one "\."
14    - (\s/$) and (~\s) are safeguards to avoid matching the end or
↳ beginning of another word
15    examples
16    -----
17    matched : M.-Madeleine Pioche de la Vergne # matched string :
↳ M.-Madeleine
18    matched : C.-A. de Ferriol # matched string : C.-A.
19    matched : J. F. # matched string : J. F.
20    matched : Jean F. # matched string : Jean F.
21    matched : Jean-F. # matched string : Jean-F.
22    matched : A M # matched string : A M
23    matched : C.-Edm.-G. # matched string : C.-Edm.-G.
24    matched : Charles-Edm.-G. # matched string : Charles-Edm.-G.
25    not matched : Anne M
26    not matched : Claude Henri blabla
27    not matched : Claude Henri
28    :param nstr: the name string used as input
29    :return: the matched string if there is a match ; None if there is
↳ no match

```

```

30     """
31     mo = re.search(r"^(|,|\s)[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûü
    ↳ üøæç]*"
32         + "\.?-[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûüüøæç]*\.(
    ↳ \s|,|$)", nstr)
    ↳ \
33     or re.search(r"^(|,|\s)[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôö
    ↳ úûüüøæç]*\."
34         + "-[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûüüøæç]*\.(
    ↳ |,|$)", nstr)
    ↳ \
35     or re.search(r"^(|,|\s)[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ]\.?
36         + "[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûüüøæç]*\.(
    ↳ \s|,
    ↳ |$)", nstr)
    ↳ \
37     or re.search(r"^(|,|\s)[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôö
    ↳ úûüüøæç]*\."
38         + "\s[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ]\.(\s|,|$)", nstr) \
39     or re.search(r"^(|,|\s)[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ]\.?
    ↳ \s[A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ]\.(\s|,|$)", nstr)
    ↳ \
40     or re.search(r"([A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ]\.){2,}", nstr) \
41     or re.search(r"^(|,|\s)([A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóô
    ↳ öúûüüøæç]*\.?-)+
42         + "([A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûüüøæç]*\.)
    ↳ (\s
    ↳ |,|$)", nstr)
    ↳ \
43     or re.search(r"^(|,|\s)([A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóô
    ↳ öúûüüøæç]*\.-)+
44         + "([A-ZÀÂÄÈÉÊËÏÔŒÛÜÝ][a-zàáâäèéêëíîïðóôöúûüüøæç]*\.(
    ↳ ?)
    ↳ (\s|,|$)",
    ↳ nstr)
45     if mo is not None:
46         return mo[0]
47     else:
48         return None

```

Code source 10 – Fonction permettant d'identifier et d'extraire un nom abrégé composé



```

1 def rgx_abvsimp(nstr):
2     """
3     try to extract a "simple" (not composed) abbreviated first name.
4     ↪ if there is no match, return None
5     pattern
6     -----
7     a capital letter (possibly followed by a certain number of
8     ↪ lowercase letters)
9     ended with a dot. (\s/$) and (^/\s) are safeguards to avoid
10    ↪ matching the beginning
11    end of another word.
12    *warning* : it can also capture parts of composed abbreviated
13    ↪ names => must be used
14    in an if-elif after trying to match a composed abbreviated name
15    examples
16    -----
17    matched : bonjour Ad. # matched string : Ad.
18    matched : J. baronne # matched string : J.
19    matched : J. F. # matched string : J.
20    matched : Jean F. # matched string : F.
21    not matched : A.-M.
22    not matched : Anne M
23    not matched : Hector
24    :param nstr: the name string used as input
25    :return:
26    """
27    mo = re.search(r"^(/\s) [A-ZÀÂÄÈÉÊËÏÔËÙÜÿ] [a-zàáâäéèëïíîïòóôóúûüø]
    ↪ æç]*\.(\\s|$|,)",
    ↪ nstr)
28    if mo is not None:
29        return mo[0]
30    else:
31        return None

```

Code source 11 – Fonction permettant de repérer et d'extraire un nom abrégé simple

```

1 def rgx_complnm(nstr):
2     """
3     try to extract a complete name from a string. if there is no
↪ match, return None
4     pattern
5     -----
6     - an uppercase letter followed by several lowercase letters ;
7     - this pattern can be repeated several times, separated by a space
↪ or "-"
8     - (\s/$) and (~\s) are safeguards to avoid matching the beginning
↪ or end of another word.
9     :param nstr: the string from which a name should be extracted
10    :return:
11    """
12    mo = re.search(r"^(~\s) [A-ZÀÂÄÈÉÊËÏÔŒÙÛÜŸ] [a-zàáâäéèëïíîïðóôöúûüø
↪ æç]+"
13    + "((\s|-) [A-ZÀÂÄÈÉÊËÏÔŒÙÛÜŸ] [a-zàáâäéèëïíîïðóôöúûüøæç]+)*($
↪ |~\s|,)",
14    ↪ nstr)
15    if mo is not None:
16        return mo[0]
17    else:
18        return None

```

Code source 12 – Fonction permettant d'identifier et d'extraire un nom complet non abrégé

```

1 <item n="264" xml:id="CAT_000327_e264">
2   <!-- ... -->
3   <name type="author">Spontini (Gaspard)</name>
4   <trait>
5     <p>célèbre compositeur, auteur de la Vestale, né en 1779, mort
6       ↪ en 1851.</p>
7   </trait>
8   <!-- ... -->
9 </item>

```

(a) La source XML-TEI

```

1 {
2   'fname': 'gaspard ',
3   'lname': 'spontini ',
4   'nobname_sts': '',
5   'status': '',
6   'dates': '1779 1851 ',
7   'function': 'writer',
8   'rebuilt': False
9 }

```

(b) Le dictionnaire d'informations structurées

```

1 gaspard spontini 1779 1851 writer
2 gaspard spontini 1851 writer
3 gaspard spontini 1779 writer
4 gaspard spontini writer
5 spontini 1779 1851 writer
6 spontini 1851 writer
7 spontini 1779 writer
8 spontini writer

```

(c) Les recherches lancées sur l'API *Wikidata*

FIGURE 3 – De la TEI à l'API : document en entrée, informations extraites et chaînes de caractères recherchées

Ci dessus (3) sont présentés l'entrée et les données produites lors de l'alignement avec *Wikidata* : en premier les données en entrée, ensuite les informations qui en sont extraites et enfin les différentes chaînes de caractères recherchées. Il est à noter que, dans la plupart des cas, seules une ou deux recherches sont faites sur l'API avant d'obtenir un résultat ; ici, huit chaînes différentes ont été construites et recherchées, ce qui représente l'intégralité de l'algorithme.

```
1 {
2   "head": {
3     "vars": [
4       "painting",
5       "paintingLabel"
6     ]
7   },
8   "results": {
9     "bindings": [
10      {
11        "painting": {
12          "type": "uri",
13          "value": "http://www.wikidata.org/entity/Q4500028"
14        },
15        "paintingLabel": {
16          "xml:lang": "en",
17          "type": "literal",
18          "value": "Circular Dance"
19        }
20      },
21      {
22        "painting": {
23          "type": "uri",
24          "value": "http://www.wikidata.org/entity/Q19882576"
25        },
26        "paintingLabel": {
27          "xml:lang": "en",
28          "type": "literal",
29          "value": "Cats (rayist percep.[tion] in rose,
30             ↪ black, and yellow)"
31        }
32      },
33      {
34        "painting": {
35          "type": "uri",
36          "value": "http://www.wikidata.org/entity/Q19883614"
```

```

37         "paintingLabel": {
38             "xml:lang": "en",
39             "type": "literal",
40             "value": "Landscape 47"
41         }
42     },
43     {
44         "painting": {
45             "type": "uri",
46             "value": "http://www.wikidata.org/entity/Q19883648"
47         },
48         "paintingLabel": {
49             "xml:lang": "en",
50             "type": "literal",
51             "value": "Rayonism Blue-Green Forest"
52         }
53     },
54     {
55         "painting": {
56             "type": "uri",
57             "value": "http://www.wikidata.org/entity/Q20200781"
58         },
59         "paintingLabel": {
60             "xml:lang": "en",
61             "type": "literal",
62             "value": "Decor for the Ballet \"Liturgie\""
63         }
64     }
65 ]
66 }
67 }

```

Code source 13 – Exemple de réponse SPARQL au format JSON pour la requête 6.1

```
1  <?xml version="1.0" ?>
2  <sparql xmlns="http://www.w3.org/2005/sparql-results#">
3
4
5      <head>
6
7
8          <variable name="painting"/>
9
10
11          <variable name="paintingLabel"/>
12
13
14      </head>
15
16
17      <results>
18
19
20          <result>
21
22
23              <binding name="painting">
24
25
26                  <uri>http://www.wikidata.org/entity/Q4500028</uri>
27
28
29              </binding>
30
31
32              <binding name="paintingLabel">
33
34
35                  <literal xml:lang="en">Circular Dance</literal>
36
37
```

```

38         </binding>
39
40
41     </result>
42
43
44     <result>
45
46
47         <binding name="painting">
48
49
50             <uri>http://www.wikidata.org/entity/Q19882576</uri>
51
52
53         </binding>
54
55
56         <binding name="paintingLabel">
57
58
59             <literal xml:lang="en">Cats (rayist percep.[tion] in
60                 ↪  rose, black, and yellow)</literal>
61
62         </binding>
63
64
65     </result>
66
67
68     <result>
69
70
71         <binding name="painting">
72
73
74             <uri>http://www.wikidata.org/entity/Q19883614</uri>

```

```
</binding>
```

```
<binding name="paintingLabel">
```

```
<literal xml:lang="en">Landscape 47</literal>
```

```
</binding>
```

```
</result>
```

```
<result>
```

```
<binding name="painting">
```

```
<uri>http://www.wikidata.org/entity/Q19883648</uri>
```

```
</binding>
```

```
<binding name="paintingLabel">
```

```
<literal xml:lang="en">Rayonism Blue-Green  
↪ Forest</literal>
```

```
</binding>
```



```

112         </result>
113
114
115     <result>
116
117
118         <binding name="painting">
119
120
121             <uri>http://www.wikidata.org/entity/Q20200781</uri>
122
123
124         </binding>
125
126
127         <binding name="paintingLabel">
128
129
130             <literal xml:lang="en">Decor for the Ballet
131                 ↪  &quot;Liturgie&quot;</literal>
132
133
134         </binding>
135
136     </result>
137
138 </results>
139
140 </sparql>

```

Code source 14 – Exemple de réponse SPARQL au format XML pour la requête 6.1

```
1 {
2   # entrées précédentes...
3   "Q108438226": {
4     "instance": [
5       "Q3305213"
6     ],
7     "instanceL": [
8       "painting"
9     ],
10    "gender": [],
11    "genderL": [],
12    "citizenship": [],
13    "citizenshipL": [],
14    "lang": [],
15    "langL": [],
16    "deathmanner": [],
17    "deathmannerL": [],
18    "birthplace": [],
19    "birthplaceL": [],
20    "deathplace": [],
21    "deathplaceL": [],
22    "residplace": [],
23    "residplaceL": [],
24    "burialplace": [],
25    "burialplaceL": [],
26    "educ": [],
27    "educL": [],
28    "religion": [],
29    "religionL": [],
30    "occupation": [],
31    "occupationL": [],
32    "award": [],
33    "awardL": [],
34    "position": [],
35    "positionL": [],
36    "member": [],
37    "memberL": [],
```

```

38     "nobility": [],
39     "nobilityL": [],
40     "workcount": [],
41     "conflictcount": [],
42     "img": [
43         "http://commons.wikimedia.org/wiki/Special:FilePath/Louis%2D
         ↳ Ode%20Frott%C3%A9.jpg"
44     ],
45     "signature": [],
46     "birth": [],
47     "death": [],
48     "title": [
49         "Louis de Frott\u00e9 (1755-1800), g\u00e9n\u00e9ral
         ↳ vend\u00e9en"
50     ],
51     "inception": [
52         "1822-01-01"
53     ],
54     "author": [],
55     "authorL": [],
56     "pub": [],
57     "pubL": [],
58     "pubplace": [],
59     "pubplaceL": [],
60     "pubdate": [],
61     "creator": [
62         "Q51077254"
63     ],
64     "creatorL": [
65         "Louise Bouteiller"
66     ],
67     "material": [
68         "Q12321255",
69         "Q296955"
70     ],
71     "materialL": [
72         "canvas",
73         "oil paint"

```

```
74     ],
75     "height": [
76         "216"
77     ],
78     "genre": [
79         "Q134307"
80     ],
81     "genreL": [
82         "portrait"
83     ],
84     "movement": [],
85     "movementL": [],
86     "creaplace": [],
87     "creaplaceL": [],
88     "viafID": [],
89     "bnfID": [],
90     "isniID": [],
91     "congressID": [],
92     "idrefID": []
93 },
94 "Q3383960": {
95     "instance": [
96         "Q5"
97     ],
98     "instanceL": [
99         "human"
100    ],
101    "gender": [
102        "Q6581097"
103    ],
104    "genderL": [
105        "male"
106    ],
107    "citizenship": [
108        "Q142"
109    ],
110    "citizenshipL": [
111        "France"
```

```

112     ],
113     "lang": [
114         "Q150"
115     ],
116     "langL": [
117         "French"
118     ],
119     "deathmanner": [],
120     "deathmannerL": [],
121     "birthplace": [
122         "Q1148468"
123     ],
124     "birthplaceL": [
125         "Vendargues"
126     ],
127     "deathplace": [],
128     "deathplaceL": [],
129     "residplace": [],
130     "residplaceL": [],
131     "burialplace": [],
132     "burialplaceL": [],
133     "educ": [],
134     "educL": [],
135     "religion": [],
136     "religionL": [],
137     "occupation": [
138         "Q47064",
139         "Q82955",
140         "Q189290"
141     ],
142     "occupationL": [
143         "military personnel",
144         "politician",
145         "military officer",
146         "officer"
147     ],
148     "award": [
149         "Q10855226",

```

```

150         "Q11593374",
151         "Q1313340"
152     ],
153     "awardL": [
154         "Grand Croix of the L\u00e9gion d'honneur",
155         "Knight of the Royal and Military Order of Saint Louis",
156         "names inscribed under the Arc de Triomphe",
157         "Grand Cross of the Legion of Honour"
158     ],
159     "position": [
160         "Q21032625",
161         "Q15904689"
162     ],
163     "positionL": [
164         "member of the Chamber of Peers",
165         "Governor of Algeria"
166     ],
167     "member": [],
168     "memberL": [],
169     "nobility": [
170         "Q165503"
171     ],
172     "nobilityL": [
173         "baron"
174     ],
175     "workcount": [],
176     "conflictcount": [
177         "1"
178     ],
179     "img": [
180         "http://commons.wikimedia.org/wiki/Special:FilePath/Berthez_
        ↪ %C3%A8ne.JPG"
181     ],
182     "signature": [],
183     "birth": [
184         "1775-03-24"
185     ],
186     "death": [

```

```

187         "1847-10-09"
188     ],
189     "title": [],
190     "inception": [],
191     "author": [],
192     "authorL": [],
193     "pub": [],
194     "pubL": [],
195     "pubplace": [],
196     "pubplaceL": [],
197     "pubdate": [],
198     "creator": [],
199     "creatorL": [],
200     "material": [],
201     "materialL": [],
202     "height": [],
203     "genre": [],
204     "genreL": [],
205     "movement": [],
206     "movementL": [],
207     "creaplace": [],
208     "creaplaceL": [],
209     "viafID": [
210         "15030322"
211     ],
212     "bnfID": [
213         "14643064j"
214     ],
215     "isniID": [
216         "0000 0000 4820 9162"
217     ],
218     "congressID": [
219         "n2005062361"
220     ],
221     "idrefID": [
222         "088623629"
223     ]
224 },

```

```
225     # entrées suivantes...
226 }
```

Code source 15 – Extrait de la base de connaissance constituée grâce à SPARQL suite à l'alignement avec des entités *Wikidata*



```

1 <listPrefixDef>
2   <prefixDef
3     ident="wd"
4     matchPattern="(Q[0-9]+)"
5     replacementPattern="https://www.wikidata.org/wiki/$1">
6     <p>
7       In the <gi>body</gi>, the <att>ref</att> attributes
8       contained in <gi>name</gi> elements are pointing to to a
9       <ref target="https://www.wikidata.org/wiki/">Wikidata</ref>
10      identifier by using the <code>wd:</code> prefix. This
11      ↪ <gi>prefixDef</gi>
12      allows to rebuilt the complete URL from a wikidata
13      ↪ identifier by
14      replacing the <code>wd:</code> prefix with:
15      <code>https://www.wikidata.org/wiki/</code>.
16    </p>
17  </prefixDef>
18 </listPrefixDef>

```

Code source 16 – Le `tei:listPrefixDef` décrivant le rôle du préfixe « wd » dans les catalogues



# Images

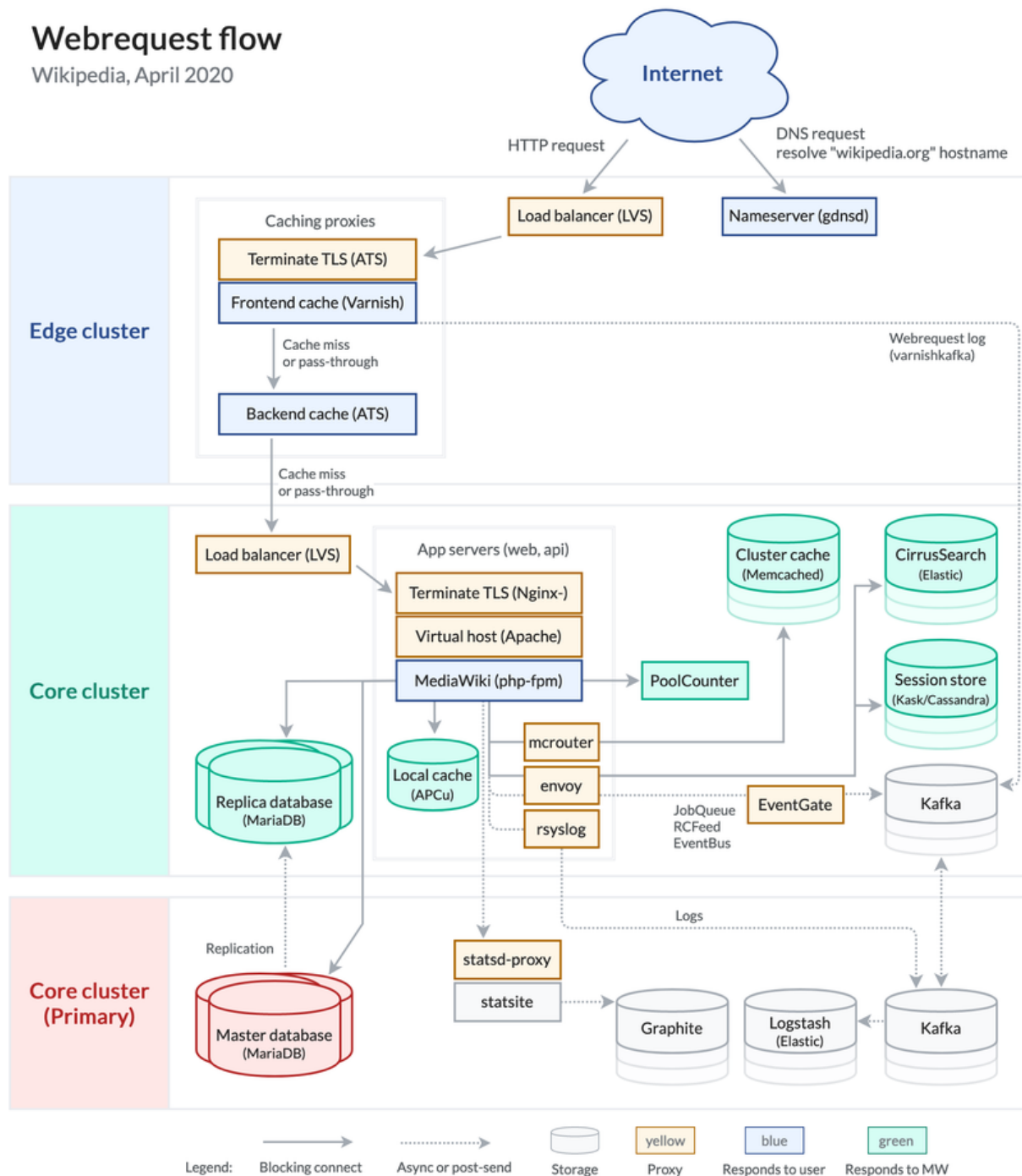


FIGURE 4 – L'architecture des bases de données de la *Wikimedia Foundation* (schéma réalisé par Timo Tijhof en 2020, disponible dans le domaine public avec la licence Creative Commons CC0.1.0 Universal Public Domain License)